

Tutorial on Interpretable Machine Learning



Wojciech Samek (Fraunhofer HHI)



Alexander Binder (SUTD)

15:00 - 15:15	Introduction & Motivation WS
---------------	------------------------------

- 15:15 16:30 Techniques for Interpretability WS
- 16:30 17:00 Coffee Break ALL
- 17:00 18:00 Applications of Interpretability WS
- 18:00 18:50 Case Study: Interpretable ML in Histopathology AB
- 18:50 19:00 Wrap-up WS



Before we start

Joint work with many people

Klaus-Robert Müller (TU Berlin) Grégoire Montavon (TU Berlin) Sebastian Lapuschkin (Fraunhofer HHI) Leila Arras (Fraunhofer HHI) Frederick Klauschen (Charite)

http://interpretable-ml.org/miccai2018tutorial/

Please ask questions at any time !



. . .

2



Tutorial on Interpretable Machine Learning

Part 1: Introduction & Motivation



MICCAI'18 Tutorial on Interpretable Machine Learning

Record Performances with ML



Jeopardy

Poker Computer games



Optical enanconversion of images on of document, a photo of a docume photo) or from subtitle text sur used as a form of information bank statements, computeris documentation. It is a comm searched, stored more com in pattern recognition, ar

4



Black Box Models

Huge volumes of data





Black Box Models



Is minimizing the error a guarantee for the model to work well in practice?



6

Why interpretability ?



MICCAI'18 Tutorial on Interpretable Machine Learning

We need interpretability in order to:

verify system understand weaknesses

8

legal aspects

learn new things from data



Why Interpretability ?

1) Verify that classifier works as expected

Wrong decisions can be costly and dangerous

"Autonomous car crashes, because it wrongly recognizes ..."



"Al medical diagnosis system misclassifies patient's disease ..."





9

2) Understand weaknesses & improve classifier





Why Interpretability ?

3) Learn new things from the learning machine

"It's not a human move. I've never seen a human play this move." (Fan Hui)



Old promise: "Learn about the human brain."





Why Interpretability ?

4) Interpretability in the sciences

Learn about the physical / biological / chemical mechanisms. (e.g. find genes linked to cancer, identify binding sites ...)







5) Compliance to legislation

European Union's new General Data Protection Regulation "right to explanation"

Retain human decision in order to assign responsibility.

"With interpretability we can ensure that ML models work in compliance to proposed legislation."



ITU/WHO Focus Group on AI4Health

Focus Group on "Artificial Intelligence for Health" established by



ITU Workshop on Artificial Intelligence for Health Geneva, Switzerland, 25 September 2018

More information about the group:

https://www.itu.int/en/ITU-T/focusgroups/ai4h



Dimensions of Interpretability

Different dimensions of "interpretability"

prediction

"Explain why a certain pattern x has been classified in a certain way f(x)."



model

"What would a pattern belonging to a certain category typically look like according to the model."





data

"Which dimensions of the data are most relevant for the task."



Dimensions of Interpretability



suboptimal or biased due to assumptions (linearity, sparsity ...)





Dimensions of Interpretability

Baehrens' Gradien	10 Sundai 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 Perte	17 urb	Kindermans'17 PatternNet
Poulin'06 Additive	Lundber Shaple Landeck	g'17 y Baze Tayl cer'13	n'13 Ma or De	ontavon'17 ep Taylor	Shrikumar'17 DeepLIFT
Zeiler'1 Decon	4 Contrib	Prop	Bach'15 LRP	Zhang Excitatior	;'16 n BP
Carua Fitted A	Sprin na'15 Gu dditi∨e	genberg'14 ided BP	Zhou'16 GAP	Selv Gra	araju'17 d-CAM

Question: Which one to choose ?





Tutorial on Interpretable Machine Learning

Part 2: Techniques of Interpretability



MICCAI'18 Tutorial on Interpretable Machine Learning

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.



Techniques of Interpretation





Model analysis



MICCAI'18 Tutorial on Interpretable Machine Learning

Approach 1: Class Prototypes

"How does a goose typically look like according to the neural network?"





Activation Maximization

- find prototypical example of a category
- find pattern maximizing activity of a neuron





Activation Maximization

- find prototypical example of a category
- find pattern maximizing activity of a neuron





Activation Maximization

- find prototypical example of a category
- find pattern maximizing activity of a neuron





Activation Maximization

Let us interpret a concept predicted by a deep neural net (e.g. a class, or a real-valued quantity):



Examples:

- Creating a class prototype: $\max_{x \in \mathcal{X}} \log p(\omega_c | x)$.
- Synthesizing an extreme case: $\max_{x \in \mathcal{X}} f(x)$.



Images from **Simonyan et al. 2013** "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps"

Observations:

- AM builds typical patterns for these classes (e.g. beaks, legs).
- Unrelated background objects are not present in the image.



Enhancing Activation Maximization

Find the input pattern that maximizes class probability.

Find the most likely input pattern for a given class.





Enhancing Activation Maximization

Images from Nguyen et al. 2016. "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks"



Observation: Connecting AM to the data distribution leads to more realistic and more interpretable images.



Application beyond Image Domain

Finding a prototype:



Question: How does a molecule with properties XYZ look like ?



Limitations of Global Interpretations

Question: Below are some images of motorbikes. What would be the best prototype to interpret the class "motorbike"?



Observations:

- Summarizing a concept or category like "motorbike" into a single image can be difficult (e.g. different views or colors).
- A good interpretation would grow as large as the diversity of the concept to interpret.



Need for Individual Explanations

Finding a prototype:



Question: How does a "motorbike" typically look like?

Individual explanation:



Question: Why is this example classified as a motorbike?



Personalized medicine: Extracting the relevant information about a medical condition for a *given* patient at a *given* time.

Each case is unique and needs its own explanation.



Personalized medicine: Extracting the relevant information about a medical condition for a *given* patient at a *given* time.

Each case is unique and needs its own explanation.

Population view: Which symptoms are most common for the disease

Both aspects can be important depending on who you are (FDA, doctor, patient).



Making Deep Neural Nets Transparent



Fraunhofer

Decision analysis



MICCAI'18 Tutorial on Interpretable Machine Learning


Sensitivity analysis: The relevance of input feature *i* is given by the squared partial derivative:

$$\mathsf{R}_i = \left(\frac{\partial f}{\partial x_i}\right)^2$$



Sensitivity analysis:







Problem: sensitivity analysis does not highlight cars

highlights parts, which (when changed) increase or decrease the prediction for "car".



Sensitivity analysis:







Problem: sensitivity analysis does not highlight cars

highlights parts, which (when changed) increase or decrease the prediction for "car".

Observation:

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

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



Shattered Gradient Problem

Input gradient (on which sensitivity analysis is based), becomes increasingly highly varying and unreliable with neural network depth.





Layer-wise Relevance Propagation (LRP)



MICCAI'18 Tutorial on Interpretable Machine Learning



Layer-wise Relevance Propagation (LRP) (Bach et al., PLOS ONE, 2015)

Explain prediction itself (not the change)









What makes this image a "rooster image" ?

Idea: Redistribute the evidence for class rooster back to image space.













Layer-wise relevance conservation

$$\sum_{i} R_{i} = \ldots = \sum_{i} R_{i}^{(l)} = \sum_{j} R_{j}^{(l+1)} = \ldots = f(x)$$



Heatmap of prediction "3"



Heatmap of prediction "9"







Explains what influences prediction "cars".

Slope decomposition $\sum_{i} R_{i} = \|\nabla_{\mathbf{x}} f\|^{2} \qquad \sum_{i} R_{i} = f(\mathbf{x})$

Explains prediction "cars" as is.

Value decomposition

More information (Montavon et al., 2017 & 2018)



Other Explanation Methods







Other Explanation Methods



Question: Which one to choose ?

(Fong & Vedaldi 2017) (Ribeiro et al., 2016) (Kindermans et al., 2017)

Deconvolution

Deconvolution Guided Backprop (Zeiler & Fergus 2014) (Springenberg et al. 2015)

Understanding the Model

Deep Visualization (Yosinski et al., 2015)

Inverting CNNs (Dosovitskiy & Brox, 2015)

Synthesis of preferred inputs (Nguyen et al. 2016)

> Network Dissection (Zhou et al. 2017)

Feature visualization (Erhan et al. 2009)

Inverting CNNs (Mahendran & Vedaldi, 2015) **RNN cell state analysis** (Karpathy et al., 2015)

Fraunhofer Heinrich Hertz Institute



Axiomatic approach to interpretability



MICCAI'18 Tutorial on Interpretable Machine Learning

First Attempt: Distance to Ground Truth





Axiomatic Approach to Interpretability

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





Axiomatic Approach to Interpretability

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).

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

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



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:





Axiomatic Approach to Interpretability

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.





Axiomatic Approach to Interpretability





Summary LRP

General Images (Bach' 15, Lapuschkin'16)



Games (Lapuschkin'18)

Faces (Lapuschkin'17)





VQA (Arras'18) there is a metallic cube ; are there any large cyan metallic objects wening it ?







Video (Anders'18)







Histopathology (Binder'18)



Text Analysis (Arras'16 &17)



Morphing (Seibold'18)



Gait Patterns (Horst'18)



fMRI (Thomas'18)



54





Summary LRP



Bag-of-words / Fisher Vector models (Bach'15, Arras'16, Lapuschkin'17, Binder'18)







Summary LRP

- 1. LRP solves the "correct" explanation problem
- 2. It has a theoretical interpretation (Deep Taylor Decomposition)
- 3. It can be applied to various data and models (not only deep nets)
- 4. It fulfills various criteria (axiomatic approach)
- 5. It is flexible (many explanation methods are special cases of LRP)
- 6. In general: LRP \neq Gradient \times Input

Tutorial Paper <u>Montavon</u> et al., "Methods for interpreting and understanding deep neural networks", Digital Signal Processing, 73:1-5, 2018

Keras Explanation Toolbox https://github.com/albermax/innvestigate



From LRP to Deep Taylor Decomposition



MICCAI'18 Tutorial on Interpretable Machine Learning

Decomposing the Correct Quantity

slope decompositionvalue decomposition
$$\sum_i R_i = \|\nabla_{\mathbf{x}} f\|^2$$
 \rightarrow $\sum_i R_i = f(\mathbf{x})$

Candidate: Taylor decomposition

$$f(\mathbf{x}) = \underbrace{f(\widetilde{\mathbf{x}})}_{0} + \sum_{i=1}^{d} \underbrace{\frac{\partial f}{\partial x_{i}}}_{R_{i}} \Big|_{\mathbf{x} = \widetilde{\mathbf{x}}} (x_{i} - \widetilde{x}_{i}) + \underbrace{O(\mathbf{x}\mathbf{x}^{\top})}_{0}$$

Achievable for linear models and deep ReLU networks without biases, by choosing:

$$\widetilde{\boldsymbol{x}} = \lim_{\varepsilon \to 0} \varepsilon \cdot \boldsymbol{x} \approx \boldsymbol{0}.$$





Why Simple Taylor doesn't work?

Two Reasons:



Root point is hard to find or too far \rightarrow includes too much information (incl. negative evidence)



Gradient shattering problem → gradient of deep nets has low informative value















Can we express R_k as a simple function of $(a_j)_j$? Can we do a Taylor decomposition of $R_k((a_i)_i)$?







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















Decompose Relevance

Taylor expansion at root point:




$R_{i\leftarrow j}$ =	(Deep Taylor generic)				
Choice of root point		$\widetilde{\pmb{a}}^{(j)} \in \mathcal{D}$	$\ \boldsymbol{a} - \widetilde{\boldsymbol{a}}^{(j)}\ $		
1. nearest root	$\widetilde{\boldsymbol{a}}^{(j)} = \boldsymbol{a} - t \cdot \boldsymbol{w}_j$		1		
2. rescaled activation	$\widetilde{a}^{(j)} = a - t \cdot a$	1			
3. rescaled excitations	$\widetilde{a}^{(j)} = a - t \cdot a \odot 1_{w_j \succ 0}$	1	1		

 $R_{i \leftarrow j} = \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j \qquad \text{(LRP-}\alpha_1 \beta_0\text{)}$



Input domain	Rule
ReLU activations $(a_j \ge 0)$	$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$
Pixel intensities $(x_i \in [l_i, h_i],$ $l_i \leq 0 \leq h_i)$	$R_{i} = \sum_{j} \frac{x_{i}w_{ij} - l_{i}w_{ij}^{+} - h_{i}w_{ij}^{-}}{\sum_{i} x_{i}w_{ij} - l_{i}w_{ij}^{+} - h_{i}w_{ij}^{-}}R_{j}$
Real values $(x_i \in \mathbb{R})$	$R_i = \sum_j \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$

Deep Taylor LRP rules [Montavon'17]

More refined rules can also be constructed to match the input data distribution [Kindermans'17]

21

68



Pooling relevance over all outgoing neurons













Tutorial on Interpretable Machine Learning

Part 3: Applications of Interpretability



MICCAI'18 Tutorial on Interpretable Machine Learning





Heinrich Hertz Institute



73

General Images (Bach' 15, Lapuschkin'16)



Games (Lapuschkin'18)

Faces (Lapuschkin'17)





VQA (Arras'18) there is a metallic cube ; are there any large cyan metallic objects wening it ?







Video (Anders'18)







Histopathology (Binder'18)



Text Analysis (Arras'16 &17)



Morphing (Seibold'18)



Gait Patterns (Horst'18)



fMRI (Thomas'18)









Bag-of-words / Fisher Vector models (Bach'15, Arras'16, Lapuschkin'17, Binder'18)







LRP & Others Evaluating Heatmap Quality



MICCAI'18 Tutorial on Interpretable Machine Learning



Can we objectively measure which heatmap is best ?

Idea: Compare selectivity (Bach'15, Samek'17):

"If input features are deemed relevant, removing them should reduce evidence at the output of the network."

Algorithm ("Pixel Flipping")

```
Sort pixels / patches by relevance
Iterate
  destroy pixel / patch
  evaluate f(x)
Measure decrease of f(x)
```

Important: Remove information in a non-specific manner (e.g. sample from uniform distribution)















LRP





























LRP:	0.722
Sensitivity:	0.691
Random:	0.523

LRP produces quantitatively better heatmaps than sensitivity analysis and random.

What about more complex datasets ?

SUN397



397 scene categories (108,754 images in total)

ILSVRC2012



1000 categories (1.2 million training images)

MIT Places



205 scene categories (2.5 millions of images)









LRP_100

Deconv. La

Deconv. l.

Random

Sensitivity

Sensitivity l

100

0.01

250

200

150

100

50

0

to random

relative

AOPC



(Samek et al. 2017)

LRP produces better heatmaps

SUN397

100

80

60

40

20

0

-20

0

20

40

60

perturbation steps

AOPC relative to random

- Sensitivity heatmaps are noisy (gradient shuttering)
- Deconvolution and sensitivity analysis solve a different problem

20

40

60

perturbation steps

80

ILSVRC2012



Same idea can be applied for other domains (e.g. text document classification)

"Pixel flipping" = "Word deleting"

Text classified as "sci.med" \rightarrow LRP identifies most relevant words.

Yes, weightlessness does feel like falling. It may feel strange at first, but the body does adjust. The feeling is not too different from that of sky diving.

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in space to try to see how to keep the number of occurances down.

(Arras et al. 2017)





⁻ word2vec / CNN model

- Conv \rightarrow ReLU \rightarrow 1-Max-Pool \rightarrow FC
- trained on 20Newsgroup Dataset

- accuracy: 80.19%

LRP better than SA

LRP distinguishes between positive and negative evidence



⁽Arras et al. 2016)







Highly efficient (e.g., 0.01 sec per VGG16 explanation) !

New Keras Toolbox available for explanation methods: https://github.com/albermax/innvestigate





Application of LRP Compare models



MICCAI'18 Tutorial on Interpretable Machine Learning

Application: Compare Classifiers

4.

Yes, weightlessness does feel like falling. It may feel strange at first, but the body does adjust. The feeling is not too different from that of sky diving.

word2vec/CNN:

Performance: 80.19%

<u>Strategy to solve the problem</u>: identify semantically meaningful words related to the topic.

BoW/SVM:

Performance: 80.10%

<u>Strategy to solve the problem</u>: identify statistical patterns, i.e., use word statistics

- >And what is the motion sickness
- >that some astronauts occasionally experience?
- It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in

Yes, weightlessness does feel like falling. It may feel strange at first, but the body does adjust. The feeling is not too different from that of sky diving.

space to try to see how to keep the number of occurances down.

- >And what is the motion sickness
- >that some astronauts occasionally experience?

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in space to try to see how to keep the number of occurances down.

(Arras et al. 2016 & 2017)



word2vec / CNN model

sci.med

symptoms (7.3), treatments (6.6), medication (6.4), osteopathy (6.3), ulcers (6.2), sciatica (6.0), hypertension (6.0), herb (5.6), doctor (5.4), physician (5.1), Therapy (5.1), antibiotics (5.1), Asthma (5.0), renal (5.0), medicines (4.9), caffeine (4.9), infection (4.9), gastrointestinal (4.8), therapy (4.8), homeopathic (4.7), medicine (4.7), allergic (4.7), dosages (4.7), esophagitis (4.7), inflammation (4.6), arrhythmias (4.6), cancer (4.6), disease (4.6), migraine (4.6), patients (4.5). BoW/SVM model

sci.med

cancer (1.4), photography (1.0), doctor (1.0), msg (0.9), disease (0.9), medical (0.8), sleep (0.8), radiologist (0.7), eye (0.7), treatment (0.7), prozac (0.7), vitamin (0.7), epilepsy (0.7), health (0.6), yeast (0.6), skin (0.6), pain (0.5), liver (0.5), physician (0.5), she (0.5), needles (0.5), dn (0.5), circumcision (0.5), syndrome (0.5), migraine (0.5), antibiotic (0.5), water (0.5), blood (0.5), fat (0.4), weight (0.4).

Words with maximum relevance

(Arras et al. 2016 & 2017)



Visual Object Classes Challenge: 2005 - 2012

	aeroplane	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	diningtable	dog	horse	motorbike	person	pottedplant	sheep	sofa	train	tymonitor
INRIA_Flat	74.8	62.5	51.2	69.4	29.2	60.4	76.3	57.6	53.1	<u>41.1</u>	54.0	42.8	76.5	62.3	84.5	35.3	41.3	50.1	77.6	49.3
INRIA_Genetic	77.5	63.6	56.1	71.9	33.1	60.6	78.0	58.8	53.5	42.6	54.9	45.8	77.5	64.0	85.9	36.3	44.7	50.6	79.2	53.2
INRIA_Larlus	62.6	54.0	32.8	47.5	17.8	46.4	69.6	44.2	44.6	26.0	38.1	34.0	66.0	55.1	77.2	13.1	29.1	36.7	62.7	43.3
MPI_BOW	58.9	46.0	31.3	59.0	16.9	40.5	67.2	40.2	44.3	28.3	31.9	34.4	63.6	53.5	75.7	22.3	26.6	35.4	60.6	40.6
PRIPUVA	48.6	20.9	21.3	17.2	6.4	14.2	45.0	31.4	27.4	12.3	14.3	23.7	30.1	13.3	62.0	10.0	12.4	13.3	26.7	26.2
QMUL_HSLS	70.6	54.8	35.7	64.5	27.8	51.1	71.4	54.0	46.6	36.6	34.4	39.9	71.5	55.4	80.6	15.8	35.8	41.5	73.1	45.5
QMUL_LSPCH	71.6	55.0	41.1	65.5	27.2	51.1	72.2	55.1	47.4	35.9	37.4	41.5	71.5	57.9	80.8	15.6	33.3	41.9	76.5	45.9
TKK	71.4	51.7	48.5	63.4	27.3	49.9	70.1	51.2	51.7	32.3	46.3	41.5	72.6	60.2	82.2	31.7	30.1	39.2	71.1	41.0
ToshCam_rdf	59.9	36.8	29.9	40.0	23.6	33.3	60.2	33.0	41.0	17.8	33.2	33.7	63.9	53.1	77.9	29.0	27.3	31.2	50.1	37.6
ToshCam_svm	54.0	27.1	30.3	35.6	17.0	22.3	58.0	34.6	38.0	19.0	27.5	32.4	48.0	40.7	78.1	23.4	21.8	28.0	45.5	31.8
Tsinghua	62.9	42.4	33.9	49.7	23.7	40.7	62.0	35.2	42.7	21.0	38.9	34.7	65.0	48.1	76.9	16.9	30.8	32.8	58.9	33.1
UVA_Bigrams	61.2	33.2	29.4	45.0	16.5	37.6	54.6	31.3	39.9	17.2	31.4	30.6	61.6	42.4	74.6	14.5	20.9	23.5	49.9	30.0
UVA_FuseAll	67.1	48.1	43.3	58.1	19.9	46.3	61.8	41.9	48.4	27.8	41.9	38.5	69.8	51.4	79.4	32.5	31.9	36.0	66.2	40.3
UVA_MCIP	66.5	47.9	41.0	58.0	16.8	44.0	61.2	40.5	48.5	27.8	41.7	37.1	66.4	50.1	78.6	31.2	32.3	31.9	66.6	40.3
UVA_SFS	66.3	49.7	43.5	60.7	18.8	44.9	64.8	41.9	46.8	24.9	42.3	33.9	71.5	53.4	80.4	29.7	31.2	31.8	67.4	43.5
UVA_WGT	59.7	33.7	34.9	44.5	22.2	32.9	55.9	36.3	36.8	20.6	25.2	34.7	65.1	40.1	74.2	26.4	26.9	25.1	50.7	29.7
XRCE	72.3	57.5	53.2	68.9	28.5	57.5	75.4	50.3	52.2	39.0	46.8	45.3	75.7	58.5	84.0	32.6	39.7	50.9	75.1	49.5



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

Test error for various classes:

(Lapuschkin et al. 2016)



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

Test error for various classes:

same performance -> same strategy ?

(Lapuschkin et al. 2016)



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

Test error for various classes:





same performance -> same strategy ?

(Lapuschkin et al. 2016)


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

Test error for various classes:

Image







same performance -> same strategy ?



Application: Compare Classifiers



'horse' images in PASCAL VOC 2007













Application: Compare Classifiers





Application: Compare Classifiers





Application of LRP Quantify Context Use



MICCAI'18 Tutorial on Interpretable Machine Learning

Application: Measure Context Use





Application: Measure Context Use

- BVLC reference model + fine tuning

- PASCAL VOC 2007





Application: Measure Context Use





Application of LRP Detect Biases & Improve Models



MICCAI'18 Tutorial on Interpretable Machine Learning

- Compare AdienceNet, CaffeNet, GoogleNet, VGG-16
- Adience dataset, 26,580 images

Age classification

	Α	С	G	V
[i]	51.4 87.0	52.1 87.9	54.3 89.1	_
[r]	51.9 87.4	52.3 88.9	53.3 89.9	-
[m]	53.6 88.4	54.3 89.7	56.2 90.7	_
[i,n]	-	51.6 87.4	56.2 90.9	53.6 88.2
[r,n]	-	52.1 87.0	57.4 91.9	-
[m,n]	-	52.8 88.3	58.5 92.6	56.5 90.0
[i,w]	-	-	-	59.7 94.2
[r,w]	_	-	-	
[m,w]	-	_	-	62.8 95.8

Gender classification

	A	С	G	V
[i]	88.1	87.4	87.9	—
[r]	88.3	87.8	88.9	_
[m]	89.0	88.8	89.7	_
[i,n]	-	89.9	91.0	92.0
[r , n]	_	90.6	91.6	
[m,n]	-	90.6	91.7	92.6
[i,w]	-	-	-	90.5
[r , w]	-	_	_	—
[m,w]	-	_	_	92.2

 $\begin{array}{l} \mathsf{A} = \mathsf{AdienceNet} \\ \mathsf{C} = \mathsf{CaffeNet} \\ \mathsf{G} = \mathsf{GoogleNet} \end{array}$

V = VGG-16

- [i] = in-place face alignment
- [r] = rotation based alignment

[m] = mixing aligned images for training

- [n] = initialization on Imagenet
- [w] = initialization on IMDB-WIKI



Gender classification



with pretraining

without pretraining

<u>Strategy to solve the problem</u>: Focus on chin / beard, eyes & hear, but without pretraining the model overfits



Age classification



Predictions

25-32 years old

<u>Strategy to solve the problem</u>: Focus on the laughing ...

60+ years old laughing speaks against 60+ (i.e., model learned that old people do not laugh)



Age classification



Predictions

25-32 years old

<u>Strategy to solve the problem:</u> Focus on the laughing ...

laughing speaks against 60+ 60+ years old (i.e., model learned that old people do not laugh) pretraining on ImageNet

pretraining on **IMDB-WIKI**





real person	fake person	real person				 1,900 images of different pretrained VGG19 mod 	nt individuals Iel
				Diff	ferent training	g methods	
				naive	one morphed	complex morphs	multiclass
			true positive	95%	90%	93%	92%
			true negative	98%	95%	95%	99%
			EER	3.1%	7.2%	6.1%	2.8%
50% (50% (genuine ima complete mo	ges, orphs			*		
	50% 10% 4 × 1	genuine ima complete m 0% one reg	ages, orphs and ion morphed	50% ge 10% cc partial r one, tw region r	enuine images, omplete morphs norphs with 10 o, three and for norphed	partial morph s, one, two, thre % morphed regi ur for two class last layer rein	s with zero, ee or four ons, classification itialized

(Seibold et al., 2018)



Semantic attack on the model	Table 4. Robustness against partial morphs.						
		left eye	right eye	nose	mouth	average	
	naive	25%	21%	14%	13%	20%	
	one morphed	81%	89%	79%	71%	80%	
	complex morphs	78%	74%	73%	54%	70%	
	multiclass	86%	93%	90%	79%	87%	
Black box adversarial attack on the model	80	-		1			



Fig. 5. Robustness against fast gradient sign attacks.



	relative amount of relevance per region								
morphed	naive					one morphed			
region	left eye	right eye	nose	mouth	left eye	right eye	nose	mouth	
left eye	0.84	0.00	0.02	0.14	0.96	0.00	0.01	0.04	
right eye	0.00	0.91	0.05	0.05	0.00	0.92	0.01	0.07	
nose	0.21	0.28	0.47	0.04	0.00	0.01	0.97	0.02	
mouth	0.34	0.27	0.04	0.35	0.17	0.12	0.04	0.68	
	с	omplex m	orphs			multiclass			
	left eye	right eye	nose	mouth	left eye	right eye	nose	mouth	
left eye	0.98	0.00	0.00	0.02	0.00	0.98	0.00	0.01	
right eye	0.00	0.92	0.00	0.08	0.98	0.00	0.02	0.00	
nose	0.02	0.03	0.92	0.02	0.01	0.10	0.19	0.70	
mouth	0.06	0.00	0.41	0.53	0.11	0.18	0.58	0.13	

(Seibold et al., 2018)







Application of LRP Learn new Representations



MICCAI'18 Tutorial on Interpretable Machine Learning

Application: Learn new Representations



(Arras et al. 2016 & 2017)



Application: Learn new Representations

2D PCA projection of document vectors





Document vector computation is <u>unsupervised</u> (given we have a classifier).

(Arras et al. 2016 & 2017)



Application of LRP Interpreting Scientific Data



MICCAI'18 Tutorial on Interpretable Machine Learning

Application: EEG Analysis

Brain-Computer Interfacing



Neural network learns that:

Left hand movement imagination leads to desynchronization over right sensorimotor cortext (and vice versa).





Application: EEG Analysis

Our neural networks are interpretable:

We can see for every trial "why" it is classified the way it is.



(Sturm et al. 2016)



Application: fMRI Analysis

Difficulty to apply deep learning to fMRI :

- high dimensional data (100 000 voxels), but only few subjects -
- results must be interpretable (key in neuroscience) -

Our approach:

- Recurrent neural networks (CNN + LSTM) for wholebrain analysis
- LRP allows to interpret the results



Dataset:

- 100 subjects from Human Connectome Project -
- N-back task (faces, places, tools and body parts) -

(Thomas et al. 2018)



0.8

1.0

Application: fMRI Analysis

	A: Group	B: Subject	C: Trial	D: TR
Body				0 0.001
Faces				
Places	0 17	05	05	0 0.0004
Tools				0 0.0015

(Thomas et al. 2018)



Application: Gait Analysis



Our approach:

- Classify & explain individual gait patterns
- Important for understanding diseases such as Parkinson





Application: Gait Analysis



Our approach:

- Classify & explain individual gait patterns
- Important for understanding diseases such as Parkinson





Application of LRP Understand Model & Obtain new Insights



MICCAI'18 Tutorial on Interpretable Machine Learning

- Fisher Vector / SVM classifier

- PASCAL VOC 2007













Motion vectors can be extracted from the compressed video -> allows very efficient analysis

- Fisher Vector / SVM classifier

- Model of Kantorov & Laptev, (CVPR'14)
- Histogram Of Flow, Motion Boundary Histogram
- HMDB51 dataset







Motion vectors can be extracted from the compressed video -> allows very efficient analysis

- Fisher Vector / SVM classifier

- Model of Kantorov & Laptev, (CVPR'14)
- Histogram Of Flow, Motion Boundary Histogram
- HMDB51 dataset







How to handle multiplicative interactions ?

 $z_j = z_g \cdot z_s$ $R_g = 0$ $R_s = R_j$ gate neuron indirectly affect relevance distribution in forward pass

Negative sentiment



Model understands negation !

(Arras et al., 2017 & 2018)



- 3-dimensional CNN (C3D)

- trained on Sports-1M

- explain predictions for 1000 videos from the test set

frame 1 frame 4 frame 7 frame 10 frame 13 frame 16



(Anders et al., 2018)





(Anders et al., 2018)





(Anders et al., 2018)




Observation: Explanations focus on the bordering

of the video, as if it wants to watch more of it.





Idea: Play video in fast forward (without retraining) and then the classification accuracy improves.



- trained on spectrograms

- spoken digits dataset (AudioMNIST)



model classifies gender based on the fundamental frequency and its immediate harmonics (see also Traunmüller & Eriksson 1995)

(Becker et al., 2018)





(Arras et al., 2018)



Sensitivity Analysis LRP 021 미큰 П

does not focus on where the ball is, but on where the ball could be in the next frame LRP shows that that model tracks the ball



Sensitivity Analysis LRP 021 미큰 П

does not focus on where the ball is, but on where the ball could be in the next frame LRP shows that that model tracks the ball





After 25 epochs



After 195 epochs







After 25 epochs



After 195 epochs























Tutorial on Interpretable Machine Learning

Part 4: Case Study: Interpretable ML in Histopathology



MICCAI'18 Tutorial on Interpretable Machine Learning

Heatmapping - a quick case study in histopathology

Talk for MICCAI workshop on Interpretable ML, 2018. Alexander Binder Joint work with F. Klauschen, S. Lapuschkin (Bach), G. Montavon, K.-R. Müller, W. Samek

ISTD Pillar, Singapore University of Technology and Design (SUTD)

September 15, 2018









Deep Neural networks and (near-)human performance

Lipnet beats humans at lip reading



Human performance in Generic classification

IM A GENET

DeepStack outplays Humans in poker



Computer outplays Humans in DOOM



human performance in low-res (!) traffic sign recognition



Mimicking art styles: https://deepart.io

Deep Learning tops human average on a constrained (!) reading comprehension task (SQuAD Dataset)

Human-like performance \neq Human-like reasoning

Adversarial attacks against deep neural networks are easy.





3.1%

Can explanation be a useful tool beyond mere curiousity?



Can explanation be a useful tool beyond mere curiousity?



Application Idea: Finding Biases in Your Training Data



Application Idea: Finding Biases in Your Training Data



Application Idea: Finding Biases in Your Training Data



Can explanation be a useful tool beyond mere curiousity?

- BoW: heatmapping for cancer evidence
- BoW: heatmapping for molecular expression evidence
- Deep Learning: heatmapping for looking for biases.

Advertisement Warning

The next slides show authors own research. Views might be positively biased ;) . Nope I have not solved all interpretability problems with it.

Why do we still talk about BoW?

- Good performance for small sample sizes (samples per class $< 10^3$).
- Stable against small changes in data augmentation / choices of negative sampling.
- Heatmapping slower compared to DNN+GPU+innvestigate

Useful for ?

Useful for ? Shows cases where heatmapping works well



No stain normalization was used here – stability. *Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological and molecular tumor profiles*, Binder et al., arxiv 2018

Useful for ? Shows cases where heatmapping works well



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

Useful for ? Shows cases where heatmapping works well



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

Useful for ? Shows cases where heatmapping works well



Useful for ? Shows cases where heatmapping fails



Too similar to dense clusters of TiLs

Useful for ? compare to:



Too similar to dense clusters of TiLs

Useful for ? Shows cases where heatmapping fails



Too similar to lymphocytes, BoW feature does not capture that their distribution is untypical for lymphos

Useful for ? Shows cases where heatmapping fails



Too similar to epithelial cells?? (patch-wise kernel similarity matrix may reveal this)

Useful for ?

Finding subtypes that are not recognized well, for example because undersampled in the training+testing set.

potential solutions:

- improved sampling
- feature engineering (BoW)
- data augmentation engineering (deep learning)

BoW: heatmapping for X

Cancer is obvious. How about molecular properties?

Example: measurement of RNA in a biopsy sample. Can we localize evidence for the expression of the corresponding protein?

Example p53, a tumor suppressor molecule

No ad-hoc localization existent.

- forward pass: predict predict concentration from HE stain as a classification problem
- backward pass: find evidence localized to pixels



BoW: heatmapping for p53



Upper column: low expression case. Lower column: High expression case. Middle row: predicted. Right row:

Immunostained groundtruth from a neighboring slice.

Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological

and molecular tumor profiles, Binder et al., arxiv 2018

31.3%
BoW: heatmapping for p53

Idea is not limited to heatmapping of evidence for cellular structures.

Accuracy is high on large subsets of patients:



32.8%

Forward pass is kernel machine over BoW feature. Backward pass to obtain scores per pixel:

- Backpropagate from output of SVM f(x) to kernel input dimension x_(d) of x
- Backpropagate from kernel input dimensions x_(d) to local features / aggregated into the BoW feature
- Backpropagate from local feature *I* to pixel *q*

Backpropagate from output to kernel input dimension. Kernel is given as:

$$f(x) = b + \sum_{i} a_{i} y_{i} k(z_{i}, x)$$
(1)
goal:
$$f(x) \approx \sum_{d} R_{d}^{(3)}(x), \text{ where}$$
(2)

 $R_d^{(3)}(x)$ is the contribution of dimension d of the test feature $x = (x_{(1)}, \ldots, x_{(D)})$ to f(x).

Backpropagate from output to kernel input dimension. Kernel is given as:

$$f(x) = b + \sum_{i} a_{i} y_{i} k(z_{i}, x)$$
(3)

goal:
$$f(x) \approx \sum_{d} R_d^{(3)}(x)$$
 (4)

In case of dimension-wise separable kernels, such as the HIK-kernel,

$$k(z, x) = \sum_{d} \min(z_{(d)}, x_{(d)})$$
(5)

$$k(z,x) = \sum_{d} k_{d}(z_{(d)}, x_{(d)})$$
(6)

$$f(x) = b + \sum_{i} a_{i} y_{i} k(z_{i}, x)$$
(7)

$$=b + \sum_{d} \sum_{i} a_{i} y_{i} k_{d}(z_{(d)}, x_{(d)})$$
(8)

$$R_d^{(3)}(x) = \frac{b}{D} + \sum_i a_i y_i k_d(z_{(d)}, x_{(d)})$$
(9)

37.5%

Backpropagate from output to kernel input dimension. Kernel is given as:

$$f(x) = b + \sum_{i} a_{i} y_{i} k(z_{i}, x)$$
(10)
goal:
$$f(x) \approx \sum_{d} R_{d}^{(3)}(x)$$
(11)

In case of differentiable kernels, such as the χ^2 -kernel,

$$k(z,x) = \exp(-\gamma \sum_{d: z_{(d)} + x_{(d)} > 0} \frac{(z_{(d)} - x_{(d)})^2}{z_{(d)} + x_{(d)}})$$
(12)

Taylor decomposition around a root $f(x_0) = 0$ is a way:

$$f(x) \approx 0 + \sum_{d} (x_{(d)} - x_{0,(d)}) \sum_{i} a_{i} y_{i} \frac{\partial k(z_{i}, x_{0})}{\partial x_{0,(d)}}$$
(13)

Backpropagate from output to kernel input dimension. Kernel is given as:

$$f(x) = b + \sum_{i} a_{i} y_{i} k(z_{i}, x)$$
(14)
goal:
$$f(x) \approx \sum_{d} R_{d}^{(3)}(x)$$
(15)

Taylor decomposition around a root $f(x_0) = 0$ is a way:

$$f(x) \approx 0 + \sum_{d} (x_{(d)} - x_{0,(d)}) \sum_{i} a_{i} y_{i} \frac{\partial k(z_{i}, x_{0})}{\partial x_{0,(d)}}$$
(16)
$$R_{d}^{(3)}(x) = (x_{(d)} - x_{0,(d)}) \sum_{i} a_{i} y_{i} \frac{\partial k(z_{i}, x_{0})}{\partial x_{0,(d)}}$$
(17)

Backpropagate from kernel input dimension to local feature. The Bow feature is a normalized sum of mappings $m_d(I)$ of local features I onto visual word dimensions:

$$x_d = c \sum_{l} m_d(l) \tag{18}$$

One example mapping is the hard assignment onto the nearest visual word among the set of all visual words $\{w_{d'}\}$:

$$m_d(l) = 1[d = \operatorname{argmin}_{d'} ||l - w_{d'}||]$$
 (19)

Note: not differentiable in *I*, so cannot use Taylor approximation again.



Backpropagate from kernel input dimension to local feature. The Bow feature is a normalized sum of mappings $m_d(I)$ of local features I onto visual word dimensions:

$$x_d = c \sum_{l} m_d(l) \tag{20}$$

Dont care how $m_d(I)$ looks like, apply special case of LRP- ϵ -rule. that would look like:

$$R^{(2)}(I) = \sum_{d} R_{d}^{(3)} \frac{m_{d}(I)}{\sum_{l'} m_{d}(I')}$$
(21)

Have to take care for those dimensions d without any weights: $\{d \mid \sum_{l} m_d(l) = 0\}$

Backpropagate from kernel input dimension to local feature. The Bow feature is a normalized sum of mappings $m_d(I)$ of local features I onto visual word dimensions:

$$x_d = c \sum_l m_d(l) \tag{22}$$

Apply special case of LRP- ϵ -rule. Have to take care for those dimensions d without any weights: $\{d \mid \sum_{l} m_{d}(l) = 0\}$:

$$Z(x) = \{ d \mid \sum_{l} m_{d}(l) = 0 \}$$

$$R^{(2)}(l) = \sum_{d \notin Z(x)} R^{(3)}_{d} \frac{m_{d}(l)}{\sum_{l'} m_{d}(l')} + \sum_{d \in Z(x)} R^{(3)}_{d} \frac{1}{\sum_{l'} 1}$$
(23)
(24)

Backpropagate from local feature to pixel

Simple idea: distribute relevance $R^{(2)}(I)$ of a local feature I equally over the support pixels q, used to compute the same local feature.

$$LF(q) = \{I \mid q \in supp(I)\} : \text{ local features which touch } q \quad (25)$$
$$R^{(1)}(q) = \sum_{I \in LF(q)} \frac{R^{(2)}(I)}{|supp(I)|}$$
(26)

Deep Learning: Heatmapping for looking for biases.



Changes in the work flow due to deep learning

- Feature engineering is dead. Learn your features from data.
- No need for tuning of features by hand.





Changes in the work flow due to deep learning

- Feature engineering is dead. Learn your features from data.
- No need for tuning of features by hand.
- Long live data augmentation engineering
- Long live hyperparameter search over grids of 37 different parameters.





Data Augmentation Engineering





brightness





contrast

color

distortion





This are 6 parameters here already for hyperparameter search here



Data Augmentation Engineering



Many hyperparameters +

- batch size
- minibatch structure
- initial learning rate
- learning rate decay
- optimizer (SGD, momentum, ADAM)

Hyperparameter search in 20-dim space

Data Augmentation Engineering

How to sample elements in minibatches ?



Histopathology in research phase: the inevitable problem of biases

Given a prediction target - example: find evidence for cancer cells.

- If a subclass is undersampled, poor performance on it cannot be detected, because it is not represented in the test set.
- extreme high variability of prediction target and of background. What are relevant subclasses?
- A relevant subclass from positive or negative labeled structures possibly undersampled, and we dont know it!

Histopathology in research phase: the inevitable problem of biases

- If a subclass is undersampled, poor performance on it cannot be detected, because it is not represented in the test set.
- Leads to a design problem:
 - how to sample positive regions for annotation? (how much of certain structures need to be sampled?)
 - how to sample negative regions for annotation?

Hypothesis: Heatmapping over large test slides may reveal undersampled structures in a qualitative manner and help in the iterative solution of the design problem. setup:

- HE stain, breast cancer
- positive annotations: positions of cancer nuclei
- negative annotations: ???



1.5:1



64.1%



65.6%



67.2%



1.5:1

68.8%



70.3%



71.9%





1.5:1

observation: different bias learned, inconsistent to ratio



setup:

- HE stain, breast cancer
- positive annotations: positions of cancer nuclei
- negative annotations: (overlapping) windows without cancer nuclei
- preprocessing: shrink image to 80%,patchsize 120, grid stride
 20
- Densenet 121, batchsize 8
- LRP- ϵ for FC layers, LRP- $\beta = 0$ for all others
- innvestigate toolbox with neuron selection index for cancer





100%



100%





100%





100%





100%
Impact of Scaling



100%





100%

Impact of Scaling

observation: 100% scaling: nuclei are too large for the fixed kernel sizes, difficult to recognize cancer, too faint heatmaps

Histopathology in application phase: heatmapping for acceptance

A classifier that simply tells a clinician: "its grade 3"

Histopathology in application phase: heatmapping for acceptance

A classifier that simply tells a clinician: "its grade 3"

Problem: in case of doubt clinician cannot validate the prediction. Classifier mistaken or clinician overlooked something?

Heatmapping allows to point the clinician to relevant regions.

Histopathology in application phase: heatmapping for acceptance

A classifier that simply tells a clinician: "its grade 3"

Problem: in case of doubt clinician cannot validate the prediction. Classifier mistaken or clinician overlooked something?

Heatmapping allows to point the clinician to relevant regions. Heatmapping allows to identify nonsensical predictions on outlier samples.

Applications III: How well does LRP scale?

DenseNet-121

Made with Keras: https://github.com/albermax/innvestigate



Thank you!

Links (LRP for LSTM for example):

http://www.heatmapping.org/

Tutorial: http://www.heatmapping.org/tutorial/

for Keras: https://github.com/albermax/innvestigate
LRPToolbox:

https://github.com/sebastian-lapuschkin/lrp_toolbox
Experimental MXnet integration:

https://github.com/sebastian-lapuschkin/lrp_toolbox/ tree/python-wip/python

Demos: https://lrpserver.hhi.fraunhofer.de/



Tutorial on Interpretable Machine Learning

Wrap-up



MICCAI'18 Tutorial on Interpretable Machine Learning

Sensitivity analysis is not the question that you would like to ask!





Take Home Messages

What works for simple models doesn't work for deep models.









vulnerable to shattered gradients











Take Home Messages

LRP works 4 all: deep models, LSTMs, kernel methods ...



LRP Explanation Framework

e people are more prone to g
The mental part is usually
y is up or down, ie: the Shu
ointed towards Earth, so the
astronauts. About 50% of t
s, and NASA has done numerou

(software, tutorials, demos, insights, applications)











 $\mathsf{LRP} \neq \mathsf{Gradient} \times \mathsf{Input}$

... except for special cases. LRP was developed among others because gradient-based methods aren't satisfying.

High flexibility: Different LRP variants, free parameters

<u>Good news</u>: No need to reimplement LRP, check our software at <u>www.heatmapping.org</u>.



Explanations can be evaluated: Pixel flipping (model agnostic) And beyond LRP and DTD

[Samek et al. IEEE TNNLS 2017]



Explanation helps to improve models



Explaining ML, Now What?



Explanation helps to find flaws in models





Getting new Insights in the Sciences

A: Group	B: Subject	C: Trial	D: TR



More information

Visit:

http://www.heatmapping.org

- Tutorials
- Software
- Online Demos



Tutorial Paper

💹 Fraunhofer

Heinrich Hertz Institute

Montavon et al., "Methods for interpreting and understanding deep neural networks", Digital Signal Processing, 73:1-5, 2018

Keras Explanation Toolbox

https://github.com/albermax/innvestigate



Tutorial / Overview Papers

G Montavon, W Samek, KR Müller. Methods for Interpreting and Understanding Deep Neural Networks. *Digital Signal Processing*, 73:1-15, 2018. **Thighly Cited Paper**

W Samek, T Wiegand, and KR Müller, Explainable Artificial Intelligence: Understanding, Visualizing and Interpreting Deep Learning Models, *ITU Journal: ICT Discoveries - Special Issue 1 - The Impact of Artificial Intelligence (AI) on Communication Networks and Services*, 1(1):39-48, 2018.

Methods Papers

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.

G Montavon, S Bach, A Binder, W Samek, KR Müller. Explaining NonLinear Classification Decisions with Deep Taylor Decomposition. *Pattern Recognition*, 65:211–222, 2017

L Arras, G Montavon, K-R Müller, W Samek. Explaining Recurrent Neural Network Predictions in Sentiment Analysis. *EMNLP'17 Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis* (WASSA), 159-168, 2017.

A Binder, G Montavon, S Lapuschkin, KR Müller, W Samek. Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers. *Artificial Neural Networks and Machine Learning – ICANN 2016*, Part II, Lecture Notes in Computer Science, Springer-Verlag, 9887:63-71, 2016.

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

Evaluation Explanations

W Samek, A Binder, G Montavon, S Lapuschkin, KR Müller. Evaluating the visualization of what a Deep Neural Network has learned. *IEEE Transactions on Neural Networks and Learning Systems*, 28(11):2660-2673, 2017.



Application to Text

L Arras, F Horn, G Montavon, KR Müller, W Samek. Explaining Predictions of Non-Linear Classifiers in NLP. *Workshop on Representation Learning for NLP,* Association for Computational Linguistics, 1-7, 2016.

L Arras, F Horn, G Montavon, KR Müller, W Samek. "What is Relevant in a Text Document?": An Interpretable Machine Learning Approach. *PLOS ONE, 12(8):e0181142*, 2017.

L Arras, G Montavon, K-R Müller, W Samek. Explaining Recurrent Neural Network Predictions in Sentiment Analysis. *EMNLP'17 Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis* (WASSA), 159-168, 2017.

L Arras, A Osman, G Montavon, KR Müller, W Samek. Evaluating and Comparing Recurrent Neural Network Explanation Methods in NLP. *arXiv*, 2018.

Application to Images & Faces

S Lapuschkin, A Binder, G Montavon, KR Müller, Wojciech Samek. Analyzing Classifiers: Fisher Vectors and Deep Neural Networks. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2912-20, 2016.

S Bach, A Binder, KR Müller, W Samek. Controlling Explanatory Heatmap Resolution and Semantics via Decomposition Depth. *IEEE International Conference on Image Processing (ICIP)*, 2271-75, 2016.

F Arbabzadeh, G Montavon, KR Müller, W Samek. Identifying Individual Facial Expressions by Deconstructing a Neural Network. *Pattern Recognition - 38th German Conference, GCPR 2016*, Lecture Notes in Computer Science, 9796:344-54, Springer International Publishing, 2016.

S Lapuschkin, A Binder, KR Müller, W Samek. Understanding and Comparing Deep Neural Networks for Age and Gender Classification. *IIEEE International Conference on Computer Vision Workshops (ICCVW)*, 1629-38, 2017.

C Seibold, W Samek, A Hilsmann, P Eisert. Accurate and Robust Neural Networks for Security Related Applications Exampled by Face Morphing Attacks. arXiv:1806.04265, 2018.



Application to Video

C Anders, G Montavon, W Samek, KR Müller. Understanding Patch-Based Learning by Explaining Predictions. *arXiv:1806.06926*, 2018.

V Srinivasan, S Lapuschkin, C Hellge, KR Müller, W Samek. Interpretable Human Action Recognition in Compressed Domain. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1692-96, 2017.

Application to Speech

S Becker, M Ackermann, S Lapuschkin, KR Müller, W Samek. Interpreting and Explaining Deep Neural Networks for Classification of Audio Signals. *arXiv:1807.03418*, 2018.

Application to Sciences

I Sturm, S Lapuschkin, W Samek, KR Müller. Interpretable Deep Neural Networks for Single-Trial EEG Classification. *Journal of Neuroscience Methods*, 274:141–145, 2016.

A Thomas, H Heekeren, KR Müller, W Samek. Interpretable LSTMs For Whole-Brain Neuroimaging Analyses. *arXiv*, 2018.

KT Schütt, F. Arbabzadah, S Chmiela, KR Müller, A Tkatchenko. Quantum-chemical insights from deep tensor neural networks. Nature communications, 8, 13890, 2017.

A Binder, M Bockmayr, M Hägele and others. Towards computational fluorescence microscopy: Machine learningbased integrated prediction of morphological and molecular tumor profiles. *arXiv:1805.11178*, 2018.

F Horst, S Lapuschkin, W Samek, KR Müller, WI Schöllhorn. What is Unique in Individual Gait Patterns? Understanding and Interpreting Deep Learning in Gait Analysis. *arXiv:1808.04308*, 2018



Software

M Alber, S Lapuschkin, P Seegerer, M Hägele, KT Schütt, G Montavon, W Samek, KR Müller, S Dähne, PJ Kindermans. iNNvestigate neural networks!. *arXiv:1808.04260*, 2018.

S Lapuschkin, A Binder, G Montavon, KR Müller, W Samek. The Layer-wise Relevance Propagation Toolbox for Artificial Neural Networks. *Journal of Machine Learning Research*, 17(114):1-5, 2016.

