Tutorial on Methods for Interpreting and Understanding Deep Neural Networks







Wojciech Samek (Fraunhofer HHI)

Grégoire Montavon Klaus-Robert Müller (TU Berlin)

(TU Berlin)

1:30 - 2:00 Part 1: Introduction

- 2:00 3:00 Part 2a: Making Deep Neural Networks Transparent
- 3:00 3:30 Break
- 3:30 4:00 Part 2b: Making Deep Neural Networks Transparent
- 4:00 5:00 Part 3: Applications & Discussion





Before we start

We thank our collaborators !



Alexander Binder (SUTD)



Sebastian Lapuschkin (Fraunhofer HHI)

Lecture notes will be online soon at:

http://www.heatmapping.org

Please ask questions at any time !



ICASSP 2017 Tutorial - W. Samek, G. Montavon & K.-R. Müller

Tutorial on Methods for Interpreting and Understanding Deep Neural Networks

W. Samek, G. Montavon, K.-R. Müller

Part 1: Introduction



Recent ML Systems achieve superhuman Performance

AlphaGo beats Go human champ



Computer out-plays humans in "doom"



💹 Fraunhofer

Heinrich Hertz Institute

Deep Net outperforms humans in image classification

IM 🔓 GENET

Autonomous search-and-rescue drones outperform humans



IBM's Watson destroys humans in jeopardy



DeepStack beats professional poker players



Deep Net beats human at recognizing traffic signs



From Data to Information

Huge volumes of data

🖉 Fraunhofer

Heinrich Hertz Institute





From Data to Information





Interpretable vs. Powerful Models ?





ICASSP 2017 Tutorial – W. Samek, G. Montavon & K.-R. Müller

Interpretable vs. Powerful Models ?



60 million parameters 650,000 neurons

We have techniques to interpret and explain such complex models !



Interpretable vs. Powerful Models ?



VS.

train interpretable model

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



Different dimensions of "interpretability"

prediction

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



"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."

model



ICASSP 2017 Tutorial - W. Samek, G. Montavon & K.-R. Müller

Why Interpretability ?

1) Verify that classifier works as expected

Wrong decisions can be costly and dangerous

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



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





Heinrich Hertz Institute



ICASSP 2017 Tutorial - W. Samek, G. Montavon & K.-R. Müller

3) Learn 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."





4) Interpretability in the sciences

Stock market analysis: "Model predicts share value with __% accuracy."

Great !!!

In medical diagnosis: "Model predicts that X will survive with probability __" What to do with this information ?



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."



Interpretability as a gateway between ML and society

- Make complex models acceptable for certain applications.
- Retain human decision in order to assign responsibility.
- "Right to explanation"

Heinrich Hertz Institut

Interpretability as powerful engineering tool

- Optimize models / architectures
- Detect flaws / biases in the data
- Gain new insights about the problem
- Make sure that ML models behave "correctly"



focus on model

Interpreting models (ensemble)

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

better understand internal representation

Explaining decisions (individual)

Fraunhofer

Heinrich Hertz Institute

- "why" does the model arrive at this particular prediction
- verify that model behaves as expected

crucial for many practical applications



In medical context

- Population view (ensemble)
 - Which symptoms are most common for the disease
 - Which drugs are most helpful for patients
- Patient's view (individual)
 - Which particular symptoms does the patient have
 - Which drugs does he need to take in order to recover

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



Interpreting models

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



 $\max_{x \in \mathcal{X}} p_{\theta}(\omega_c \,|\, x) + \lambda \Omega(x)$



Interpreting models

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



ICASSP 2017 Tutorial - W. Samek, G. Montavon & K.-R. Müller

Interpreting models

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





ICASSP 2017 Tutorial - W. Samek, G. Montavon & K.-R. Müller

Explaining decisions

- "why" does the model arrive at a certain prediction
- verify that model behaves as expected





Explaining decisions

Heinrich Hertz Institute

- "why" does the model arrive at a certain prediction
- verify that model behaves as expected



Sensitivity Analysis (Simonyan et al. 2014)





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







Sensitivity Analysis:

"what makes this image less / more 'scooter' ?"

LRP / Taylor Decomposition:

"what makes this image 'scooter' at all ?"



More to come

