

iNNvestigate neural networks!

Maximilian Alber

MAXIMILIAN.ALBER@TU-BERLIN.DE

*Technische Universität Berlin, Machine Learning Group
10623 Berlin, Germany*

Sebastian Lapuschkin

SEBASTIAN.LAPUSCHKIN@HHI.FRAUNHOFER.DE

*Fraunhofer Heinrich Hertz Institute, Video Coding and Analytics
10587 Berlin, Germany*

Philipp Seegerer

PHILIPP.SEEGERER@TU-BERLIN.DE

Miriam Hägele

HAEGELE@TU-BERLIN.DE

Kristof T. Schütt

KRISTOF.SCHUETT@TU-BERLIN.DE

Grégoire Montavon

GREGOIRE.MONTAVON@TU-BERLIN.DE

*Technische Universität Berlin, Machine Learning Group
10623 Berlin, Germany*

Wojciech Samek

WOJCIECH.SAMEK@HHI.FRAUNHOFER.DE

*Fraunhofer Heinrich Hertz Institute, Video Coding and Analytics
10587 Berlin, Germany*

Klaus-Robert Müller

KLAUS-ROBERT.MUELLER@TU-BERLIN.DE

*Technische Universität Berlin, Machine Learning Group
10623 Berlin, Germany
Korea University, Department of Brain and Cognitive Engineering
Seoul 02841, Korea
Max Planck Institute for Informatics
66123 Saarbrücken, Germany*

Sven Dähne

SVEN.DAEHNE@TU-BERLIN.DE

Pieter-Jan Kindermans

P.KINDERMANS@TU-BERLIN.DE

*Technische Universität Berlin, Machine Learning Group
10623 Berlin, Germany*

Abstract

In recent years, deep neural networks have revolutionized many application domains of machine learning and are key components of many critical decision or predictive processes. Therefore, it is crucial that domain specialists can understand and analyze actions and predictions, even of the most complex neural network architectures. Despite these arguments neural networks are often treated as black boxes. In the attempt to alleviate this shortcoming many analysis methods were proposed, yet the lack of reference implementations often makes a systematic comparison between the methods a major effort. The presented library *iNNvestigate* addresses this by providing a common interface and out-of-the-box implementation for many analysis methods, including the reference implementation for PatternNet and PatternAttribution as well as for LRP-methods. To demonstrate the versatility of *iNNvestigate*, we provide an analysis of image classifications for variety of state-of-the-art neural network architectures.

1. Introduction

In recent years deep neural networks have revolutionized many domains, e.g., image recognition, speech recognition, speech synthesis, and knowledge discovery (Krizhevsky et al., 2012; LeCun et al., 2012; Schmidhuber, 2015; LeCun et al., 2015; Van Den Oord et al., 2016). Due to their ability to naturally learn from structured data and exhibit superior performance, they are increasingly used in practical applications and critical decision processes, such as novel knowledge discovery techniques, autonomous driving or medical image analysis. To fully leverage their potential it is essential that users can *comprehend and analyze* these processes. E.g., in neural architecture (Zoph et al., 2017) or chemical compound searches (Montavon et al., 2013; Schütt et al., 2017) it would be extremely useful to know which properties help a neural network to choose appropriate candidates. Furthermore for some applications understanding the decision process might be a legal requirement.

Despite these arguments neural networks are often treated as black boxes, because their complex internal workings and the basis for their predictions are not fully understood. In the attempt to alleviate this shortcoming several methods were proposed, e.g., Saliency Map (Baehrens et al., 2010; Simonyan et al., 2013), SmoothGrad (Smilkov et al., 2017), IntegratedGradients (Sundararajan et al., 2017), Deconvnet (Zeiler and Fergus, 2014), GuidedBackprop (Springenberg et al., 2015), PatternNet and PatternAttribution (Kindermans et al., 2018), LRP (Bach et al., 2015; Lapuschkin et al., 2016a,b; Montavon et al., 2018), and DeepTaylor (Montavon et al., 2017). Theoretically it is not clear which method solves the stated problems best, therefore an empirical comparison is required (Samek et al., 2017; Kindermans et al., 2017). In order to evaluate these methods, we present *iNNvestigate* which provides a common interface to a variety of analysis methods.

In particular, *iNNvestigate* contributes:

- A common interface for a growing number of analysis methods that is applicable to a broad class of neural networks. With this instantiating a method is as uncomplicated as passing a trained neural network to it and allows for easy qualitative comparisons of methods. For quantitative evaluations of (image) classification task we further provide an implementation of the method “perturbation analysis” (Samek et al., 2017).
- Support of all methods listed above—this includes the first reference implementation for PatternNet and PatternAttribution and an extended implementation for LRP—and an open source repository for further contributions.
- A clean and modular implementation, casting each analysis in terms of layer-wise forward and backward computations. This limits code redundancy, takes advantage of automatic differentiation, and eases future integration of new methods.

iNNvestigate is available at repository: <https://github.com/albermax/innvestigate>. It can be simply installed as Python package and contains documentation for code and applications. To demonstrate the versatility of *iNNvestigate* we provide examples for the analysis of image classifications for a variety of state-of-the-art neural networks.

Terminology The different methods pose different assumption to tasks and are designed for different objectives, yet they are related to “explaining” or “interpreting” neural networks (see Montavon et al. (2018)). We actively refrain from using this terminology in order

to prevent misunderstandings between the design choices of the algorithms and the implicit assumption these terms bring along. Therefore we will solely use the neutral term *analyzing* and leave any interpretation to the user.

2. Library

Interface The main feature is a common interface to several analysis methods. The workflow is as simple as passing a Keras neural network model to instantiate an analyzer object for a desired algorithm. Then, if needed, the analyzer will be fitted to the data and eventually be used to analyze the model’s predictions. The corresponding Python code is:

```
1 import innvestigate
2 model = create_a_keras_model()
3 analyzer = innvestigate.create_analyzer("analyzer_name", model)
4 analyzer.fit(X_train) # if needed
5 analysis = analyzer.analyze(X_test)
```

Implemented methods At publication time the following algorithms are supported: Gradient Saliency Map, SmoothGrad, IntegratedGradients, Deconvnet, GuidedBackprop, PatternNet and PatternAttribution, DeepTaylor, and LRP including LRP-Z, -Epsilon, -AlphaBeta. In contrast, current related work (Kotikalapudi et al., 2017; Ancona et al., 2018) is limited to gradient-based methods. We intend to further extend this selection and invite the community to contribute implementations as new methods emerge.

Documentation The library’s documentation contains several introductory scripts and example applications. We demonstrate how the analyses can be applied to the following state-of-the-art models: VGG16 and VGG19 (Simonyan and Zisserman, 2014), InceptionV3 (Szegedy et al., 2016), ResNet50 (He et al., 2016), InceptionResNetV2 (Szegedy et al., 2017), DenseNet (Huang et al., 2017), NASNet mobile, and NASNet large (Zoph et al., 2017). Figure 1 shows the result of each analysis on a subset of these networks.

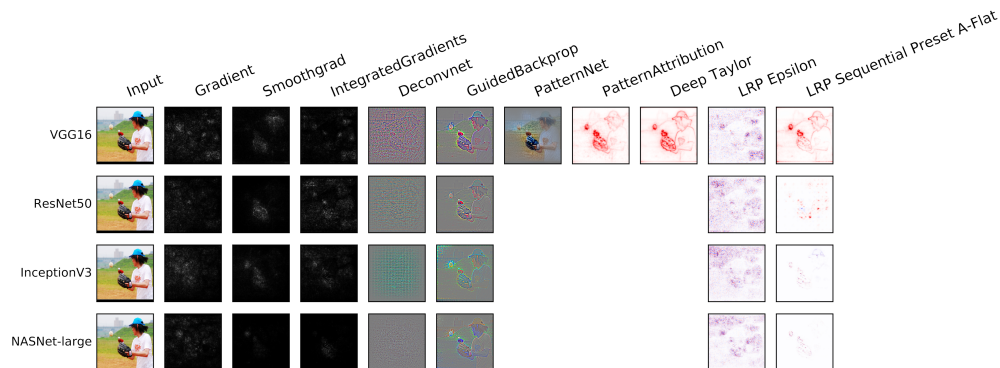


Figure 1: Result of methods applied to various neural networks (blank, if a method does not support a network’s architecture yet).

2.1 Details

Modular implementation All of the methods have in common that they perform a back-propagation from the model outputs to the inputs. The core of *iNNvestigate* is a set of base classes and functions that is designed to allow for rapid and easy development of such algorithms. The developer only needs to implement specific changes to the base algorithm and the library will take care of the complex and error-prone handling of the propagation along the graph structure. Further details can be found in the repositories documentation.

Another advantage of the modular design is that one can extend any analyzer with a given set of wrappers. One application of this is the smoothing of the analysis results by adding Gaussian noise to the copies of the input and averaging the outcome. E.g., SmoothGrad is realized in this way by combining a smoothing wrapper with a gradient analyzer.

Training PatternNet and PatternAttribution (Kindermans et al., 2018) are two novel approaches that condition their analysis on the data distribution. This is done by identifying the signal and noise direction for each neuron of a neural network. Our software scales favorably, e.g., one can train required patterns for the methods on large datasets like Imagenet (Deng et al., 2009) in less than an hour using one GPU. We present the first reference implementation of these methods.

Quantitative evaluation Often analysis methods for neural networks are compared by qualitative (visual) inspection of the result. This can lead to subjective evaluations and one approach to create a more objective and quantitative comparison of analysis algorithms is the method “perturbation analysis” (Samek et al., 2017, also known as “PixelFlipping”). The intuition behind this method is that perturbing regions which are recognized as important for the classification task by the analyzing method, will impact the classification most. This allows to assess which analysis method best identifies regions that matter for a specific task and neural network. *iNNvestigate* contains an implementation of this method.

Installation & license *iNNvestigate* is published as open-source software under the MIT-license and can be downloaded from: <https://github.com/albermax/innvestigate>. It is build as a Python 2 or 3 application on top of the popular and established Keras (Chollet et al., 2015) framework. This allows to use the library on various platforms and devices like CPUs and GPUs. At the time of publication only the TensorFlow (Abadi et al., 2016) Keras-backend is supported. The library can be simply installed as Python package.

3. Conclusion

We have presented *iNNvestigate*, a library that makes it easier to analyze neural networks’ predictions and to compare different analysis methods. This is done by providing a common interface and implementations for many analysis methods as well as making tools for training and comparing methods available. In particular it contains reference implementations for many methods (PatternNet, PatternAttribution, LRP) and example application for a large number of state-of-the-art applications. We expect that this library will support the field of analyzing machine learning and facilitate research using neural networks in domains such as drug design or medical image analysis.

Acknowledgments

Correspondence to MA, SL, KRM, WS and PJK. This work was supported by the Federal Ministry of Education and Research (BMBF) for the Berlin Big Data Center BBDC (01IS14013A). Additional support was provided by the BK21 program funded by Korean National Research Foundation grant (No. 2012-005741) and the Institute for Information & Communications Technology Promotion (IITP) grant funded by the Korea government (no. 2017-0-00451, No. 2017-0-01779).

References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- Marco Ancona, Enea Ceolini, Cengiz ztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Sy21R9JAW>.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE*, 10(7):e0130140, 2015.
- David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(Jun):1803–1831, 2010.
- François Chollet et al. Keras. <https://keras.io>, 2015.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4700–4708, 2017.
- Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un)reliability of saliency methods. *NIPS 2017 Workshop - Interpreting, Explaining and Visualizing Deep Learning - Now what?*, 2017.
- Pieter-Jan Kindermans, Kristof T. Schtt, Maximilian Alber, Klaus-Robert Mller, Dumitru Erhan, Been Kim, and Sven Dhne. Learning how to explain neural networks: PatternNet and PatternAttribution. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Hkn7CBaTW>.
- Raghavendra Kotikalapudi et al. keras-vis. <https://github.com/raghakot/keras-vis>, 2017.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.

- Sebastian Lapuschkin, Alexander Binder, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Analyzing classifiers: Fisher vectors and deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2912–2920, 2016a.
- Sebastian Lapuschkin, Alexander Binder, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. The layer-wise relevance propagation toolbox for artificial neural networks. *Journal of Machine Learning Research*, 17(114):1–5, 2016b.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.
- Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient back-prop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 2012.
- Grégoire Montavon, Matthias Rupp, Vivekanand Gobre, Alvaro Vazquez-Mayagoitia, Katja Hansen, Alexandre Tkatchenko, Klaus-Robert Müller, and O Anatole Von Lilienfeld. Machine learning of molecular electronic properties in chemical compound space. *New Journal of Physics*, 15(9):095003, 2013.
- Grégoire Montavon, Sebastian Bach, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
- Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.
- Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert 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.
- Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61: 85–117, 2015.
- Kristof T Schütt, Farhad Arbabzadah, Stefan Chmiela, Klaus Robert Müller, and Alexandre Tkatchenko. Quantum-chemical insights from deep tensor neural networks. *Nature Communications*, 8:13890, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*, abs/1312.6034, 2013. URL <http://arxiv.org/abs/1312.6034>.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: Removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Jost T Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. In *ICLR (workshop track)*, 2015.

- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3319–3328, 2017.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence (AAAI)*, pages 4278–4284, 2017.
- Aäron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 818–833. Springer, 2014.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. *arXiv preprint arXiv:1707.07012*, 2017.