

Introspection of Neural Networks Deep Learning School ,Basic Concepts'







Credits: Yasuhiro Kubota, tf-keras-vis



Outline of this course

Block 1: Introduction & Visualisation

- 1) Overview of concepts
- 2) Feature Visualisation as a tool

Block 2: Attention & Reaction

- 1) Class Activation Maps (CAM)
- 2) Saliency
- 3) From sensitivity to attribution

Block 3: Further concepts for introspection

Wrap-up & summary

In total: ca. 1/3 lecture, 2/3 hands-on



The tutorials provide examples of neural network introspection written in Python, using the Keras* library and TensorFlow tensor ordering convention. We will also use additional toolkits such as tf-keras-vis.**

*Keras provides a high level API to create deep neural networks and train them using numerical tensor libraries (backends) such as TensorFlow, CNTK or Theano.



** https://keisen.github.io/tf-keras-vis-docs/index.html#



Credit: A. Boucaud



For our hands-on tutorial on

Introspection of neural networks:

You can start a binder image here: <u>https://mybinder.org/v2/gh/csheneka/introspection-tutorial/HEAD</u> Please start with: <u>visualisation_3D-21cmPIE-Net.ipynb</u>

OR via google Colab here: <u>https://colab.research.google.com</u> for path <u>https://github.com/csheneka/introspection-tutorial</u>

Extra data: <u>https://cloud.hs.uni-hamburg.de/s/gYWW4PQG2b57XK3</u> (larger file version)



Network models are data driven.

Need introspection to:

- Evaluate training data
- Debug, detect 'anomalies'
- Understand and verify your model
- Understand predictions



Hans am Tretbrett (1909)





Lapuschkin et al. 2019, arXiv: 1902.10178



Network models are data driven.

Need introspection to:

- Evaluate training data
- Debug, detect 'anomalies'
- Understand and verify your model
- Understand predictions

Famous examples: 'Clever Hans' effect (wrong cue, 'overfitting') Over-confidence (e.g. miss-classifying)



Lapuschkin et al. 2019, arXiv: 1902.10178









We do know: architecture (model), best-fitting parameters (weights/bias/filter) + how these change during training + reaction to data \longrightarrow analyse





To learn about the model and its reaction to data, one can:

- Visualise feature space [Block 1]
- Analyse attention & reaction (CAM, saliency) [Block 2]
- Be model-agnostic (SHAP, clustering) [Block 3]
- Many more: Put errors/probabilities, regularise, unsupervised ... [Block 3]
 ... very active field of development!

To learn about the model and its reaction to data, one can:

- Visualise feature space [Block 1]
- Analyse attention & reaction (CAM, saliency) [Block 2]
- Be model-agnostic (e.g. SHAP) [Block 3]
- Many more: Put errors/probabilities, regularise, unsupervised ... [Block 3]
 ... active field of development!
 - Visualise structure of weights (dense), filters (convolutional)
 Example we will learn about with hands-on exercises:
 Filter and activation of CNN, (activation maximisation)



Reminder: Convolutional Neural Networks

Neuron response to visual stimuli



11./12.8.2022, C. Heneka, DL Basic Concepts: Introspection

size color

bka



Feature Visualisation



How do weights/biases/filters look like? How do networks 'see'? i.e. how do they learn representations? Interpretability:

Are shapes, that we expect to be important, there?



Feature Visualisation layer by layer



Visualisation of feature maps, or activation

Input as treated by model (filters)





Transposed Convolution: Feature 'Deconvolution'

Useful for:

- Segmentation, localisation, super-resolution through upsampling
- Relating input patterns to feature maps See Zeiler et al., 2013

Goal: Map intermediate features to input shape Transposed Convolution consists of <u>unpooling</u> to create a larger, but sparse, activation map, that then is via <u>'fractional' convolution</u> connects for every input activation multiple outputs.





Transposed Convolution: Visualisation

For visualisation we are interested in mapping what caused the highest activation in feature space:

- Model: transposed trained model (weights, filters)
- Upsample only highest activation
- Up to input space dimension

+ Further examples arXiv:1804.11191



Suppression of background, specific representation, complexity

Noh, Hong, Han 2015 arXiv: 1505.04366





versität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

Useful to visualise both dense and convolutional networks parts.

An input pattern image x^* is synthesized that maximises the activation of a feature, i.e. that finds preferred input for a neuron, filter, or layer (deep dream):

 $x^* = \operatorname{argmax} a_{i,l}(\theta, x)$

for a pretrained model with parameters θ , and input x.

Steps:

(1) Start with input noise image x0

(2) Calculate activation gradients via backpropagation at fixed θ

(3) Iteratively maximise the activation, until x* $x \leftarrow x + \eta \cdot \frac{\partial a_{i,l}(\theta, x)}{2}$







- Introspection important to verify and debug your model
- ✓ Deep Learning is data-driven: need to understand data representation
- ✓ Goal is interpretable / explainable machine learning
- ✓ Introspection is possible via visualisation of trained properties
 - Visualisation of filter and feature maps
 - Transposed Convolution upsamples and identifies 'important' input shapes
 - Activation Maximisation finds shapes that a specific neuron/filter/layer is sensitive to
 - and more: network inversion, network dissection, pruning, ...
 - → We can track feature hierarchy, representation learning in a probabilistic manner.



Hands-on background: Visualisation of an 'Astrophysical' CNN

What is the cosmology at high redshifts? 'gap' between CMB and galaxy surveys What properties do the very first stars and galaxies have?





Hands-on background: Visualisation of an 'Astrophysical' CNN



Moving from 2D to full 3D convolution



Best-performing: simple Conv3D architecture





For our hands-on tutorial on

Introspection of neural networks:

You can start a binder image here: <u>https://mybinder.org/v2/gh/csheneka/introspection-tutorial/HEAD</u> Please start with: <u>visualisation_3D-21cmPIE-Net.ipynb</u>

OR via google Colab here: <u>https://colab.research.google.com</u> for path <u>https://github.com/csheneka/introspection-tutorial</u>

Extra data: <u>https://cloud.hs.uni-hamburg.de/s/gYWW4PQG2b57XK3</u> (larger file version)



Visualisation of an 'Astrophysical' CNN

Hands-on: Visualisation of filters layer by layer

A) Raw filters (spatial dim): example from Gillet et al. 2019, arXiv:1805.02699

Question: Could we do something with the filters, knowing the variation in 2D (spatial) is very different from the +1D temporal dimension?





Visualisation of an 'Astrophysical' CNN

Hands-on: Visualisation of filters layer by layer

A) Raw filters (spatial dim): example from Gillet et al. 2019, arXiv:1805.02699

Question: Could we do something with the filters, knowing the variation in 2D (spatial) is very different from the +1D temporal dimension?



B) Average in redshift (temporal) direction: Note the filters for 'slow transitions' and fluctuations (example our code)





Hands-on: Visualisation of feature maps, or activation, layer by layer





Hands-on: Visualisation of feature maps, or activation, layer by layer







Hands-on: Visualisation of feature maps, or activation, layer by layer





Part 2: Attention & Reaction

To learn about the model and its reaction to data, one can:

- Visualise feature space [Block 1]
- Analyse attention & reaction (CAM, saliency) [Block 2]
- Be Model-agnostic (e.g. SHAP) [Block 3]
- Many more: Put errors/probabilities, regularise, unsupervised ... [Block 3]
 ... active field of development!
- Analyse the attention over input for the <u>prediction</u>
 i.e. <u>sensitivity</u> or <u>attribution</u> of the prediction to input

Example we will learn about with hands-on exercises: CAM, saliency map



Lapuschkin et al. 2019, arXiv: 1902.10178



Discriminative localisation with CAM

CAM = Class Activation Map

Specific architecture: GAP* output of the last convolutional layer before dense layers.

* Global Average Pooling



Credits: Yasuhiro Kubota, tf-keras-vis

CAM = (Weighted) sum of GAP weights and activations of last CNN layer, upsampled via interpolation to input data shape



Discriminative localisation with CAM

CAM = (Weighted) sum of GAP weights and activations of last CNN layer, upsampled via interpolation to input data shape





Credits: Yasuhiro Kubota, tf-keras-vis

$w_1 \times f_1 + w_2 \times f_2 + \dots + w_N \times f_N$ + up-sample = CAM



Saliency Maps

Gradient-based approach to localisation "Most important input features (pixels) cause the highest gradient" Flexible wrt architecture Map ranks importance



Credits: Yasuhiro Kubota, tf-keras-vis

Principle: Backpropagation to input (training uses backpropagation to 1st layer)

(1) For fixed model parameters θ , obtain for input space x the NN prediction $p(x,\theta)$

(2) Calculate the gradient map (same shape as input):

$$S = \frac{\partial p(x,\theta)}{\partial x} \bigg|_{x=x_0}$$



From sensitivity to attribution

Improvements of sensitivity,

e.g. to reduce noise and achieve better localisation:

- SmoothGrad: add noise to input for ensemble of maps, average map of attention
- CAM: different algorithms for calculation of class activation (GradCAM vs. ScoreCAM, ++)

Integrated Gradients (Sundararajan, Taly, Yan 2017 arXiv:1703.01365)

Attribution methods should satisfy Sensitivity and Implementation Invariance





Credits: Yasuhiro Kubota, tf-keras-vis



Attribution methods should satisfy Sensitivity and Implementation Invariance

Attribution = How predictions are formed

Example methods: Integrated Gradients, <u>Layer-wise relevance propagation</u> (LRP), Pixel perturbation analysis

General idea:

Rank input by its relevance for the prediction, Backward mapping & conservation principle $\sum R_i^{(l)} = \sum R_j^{(l+1)}$ I = layer, neurons I, j

The Principle is based on Taylor decomposition & can in principle be applied to any architecture.

- Perform a standard forward pass for predictions
- Then: Propagate the model output backward, layer-wise, using decomposition as:

$$R_{i}^{(l)} = \sum_{j} \frac{z_{ij}}{\sum_{i'} z_{i'j} + \epsilon \operatorname{sign}(\sum_{i'} z_{i'j})} R_{j}^{(l+1)}$$





Samek et al. 2015 arXiv:1509.06321

- e Random Segmentation Relevance
- Final product: Relevance heatmap of input shape

- ✓ Introspection is possible via sensitivity analyses and localisation
 - Class Activation Maps upsample last CNN-layer to input data shape
 - Saliency is a gradient-based approach to rank importance of input pixels
- Introspection is possible by analysing how predictions are formed attribution
 Importance of sensitivity and implementation invariance
 - Example of layer-wise relevance propagation for input-shaped relevance heatmaps
- Go beyond visualization with both sensitivity and attribution to understand the input – model – prediction relations

This exercise will create and analyse saliency maps for the trained 3D-21cmPIE-Net here:

<u>https://github.com/stef-neu/3D-21cmPIE-Net/blob/main/paper_plots/saliency_maps.py</u> This uses tf_keras_vis.saliency

We will also briefly compare saliency maps with CAM generated as well with tf_keras_vis

(a github with a pedagogical walk-through as a jupyter-notebook of this will follow)



Hands-on: Attention of an 'Astrophysical' CNN

Hands-on: Creation of saliency maps and their analysis

A) Simulations (no noise)

Question: What do you notice, which areas spatially and temporally are important? How does attention shift with the inclusion of noise?





Hands-on: Attention of an 'Astrophysical' CNN





To learn about the model and its reaction to data, one can:

- Visualise feature space [Block 1]
- Analyse attention & reaction (CAM, saliency) [Block 2]
- Be Model-agnostic (e.g. SHAP) [Block 3]
- Many more: Put errors/probabilities, regularise, unsupervised ... [Block 3]
 ... active field of development!
- Only take into account relation input to output

Example we will learn about with hands-on exercises: SHAP



Only take into account relation input & predictions for attribution.

Examples:

- SHAP values
- LIME
- QII

SHAP = SHapley Additive exPlanations

Idea:

Show importance of feature on predictions (without any information on predictive quality).



SHAP = SHapley Additive exPlanations Approach from Game Theory

Idea:

Assign to each input feature a value. A larger value indicates higher importance for the output prediction.

Do so via optimal credit allocation of possible 'coalitions'.





To learn about the model and its reaction to data, one can:

- Visualise feature space [Block 1]
- Analyse attention & reaction (CAM, saliency) [Block 2]
- Model-agnostic (e.g. SHAP) [Block 3]
- Many more: Put errors/probabilities, regularise, unsupervised ... [Block 3]

... active field of research & development!



Some concepts interesting for introspection worth to be highlighted:

- Bottlenecks (e.g. in variational autoencoders): Disentanglement & physical interpretation
- Generative methods for reconstruction errors
- Regulatory terms & reaction
- <u>Uncertainty estimation</u> 'I do not know' 'Not sure' error on predictions



Part 3: Further concepts

Some concepts interesting for introspection worth to be highlighted:

- Bottlenecks (e.g. in variational autoencoders): Disentanglement & physical interpretation
- Generative methods for reconstruction errors
- Regulatory terms & reaction
- <u>Uncertainty estimation</u> 'I do not know' 'Not sure' error on predictions





Hands-on Part 3: Further concepts

Very brief exercise planned using SHAP (https://shap.readthedocs.io/en/latest/index.html)

Plus bonus if time permits: Playing e.g. with deep dream



Summary of this lecture

- We learned that <u>Introspection</u> important to verify and debug your model, together with understanding your data and predictions.
- <u>Visualisation</u> helps to learn more about the trained model; we can track to hierarchy of features in network models.
- Sensitivity and attribution relate predictions to input, enable discriminative localisation.
- 4) A zoo of methods waits to be used and developed for introspection!



