

---

# iNNvestigate Documentation

*Release not set*

**Maximilian Alber**

**Mar 19, 2019**



---

## Contents:

---

|   |           |
|---|-----------|
| <b>1 API reference</b>                  | <b>3</b>  |
| 1.1 innvestigate . . . . .              | 3         |
| 1.2 innvestigate.analyzer . . . . .     | 3         |
| 1.3 innvestigate.applications . . . . . | 11        |
| 1.4 innvestigate.tools . . . . .        | 12        |
| 1.5 innvestigate.utils . . . . .        | 15        |
| <b>2 Indices and tables</b>             | <b>19</b> |
| <b>Python Module Index</b>              | <b>21</b> |



This is the API documentation for the iNNvestigate library. It contains a references for the most important functions and classes of the library. For a usage guide and more information about the interplay of the componentes please visit the [GitHub repo](#).



# CHAPTER 1

---

## API reference

---

### 1.1 innvestigate

`innvestigate.analyzer.create_analyzer(name, model, **kwargs)`

Instantiates the analyzer with the name ‘name’

This convenience function takes an analyzer name creates the respective analyzer.

Alternatively analyzers can be created directly by instantiating the respective classes.

#### Parameters

- **name** – Name of the analyzer.
- **model** – The model to analyze, passed to the analyzer’s `__init__`.
- **kwargs** – Additional parameters for the analyzer’s .

**Returns** An instance of the chosen analyzer.

**Raises** `KeyError` – If there is no analyzer with the passed name.

### 1.2 innvestigate.analyzer

#### 1.2.1 Interface and base code

`exception innvestigate.analyzer.base.NotAnalyzableModelError`

Indicates that the model cannot be analyzed by an analyzer.

`class innvestigate.analyzer.base.AnalyzerBase(model, disable_model_checks=False)`

The basic interface of an iNVESTIGATE analyzer.

This class defines the basic interface for analyzers:

```
>>> model = create_keras_model()
>>> a = Analyzer(model)
>>> a.fit(X_train) # If analyzer needs training.
>>> analysis = a.analyze(X_test)
>>>
>>> state = a.save()
>>> a_new = A.load(*state)
>>> analysis = a_new.analyze(X_test)
```

## Parameters

- **model** – A Keras model.
- **disable\_model\_checks** – Do not execute model checks that enforce compatibility of analyzer and model.

---

**Note:** To develop a new analyzer derive from [AnalyzerNetworkBase](#).

---

### **analyze** (*X*)

Analyze the behavior of model on input *X*.

**Parameters** **X** – Input as expected by model.

#### **fit** (\*args, \*\*kwargs)

Stub that eats arguments. If an analyzer needs training include [TrainerMixin](#).

**Parameters** **disable\_no\_training\_warning** – Do not warn if this function is called despite no training is needed.

#### **fit\_generator** (\*args, \*\*kwargs)

Stub that eats arguments. If an analyzer needs training include [TrainerMixin](#).

**Parameters** **disable\_no\_training\_warning** – Do not warn if this function is called despite no training is needed.

#### **static load** (*class\_name*, *state*)

Resembles an analyzer from the state created by `analyzer.save()`.

**Parameters**

- **class\_name** – The analyzer's class name.
- **state** – The analyzer's state.

#### **static load\_npz** (*fname*)

Resembles an analyzer from the file created by `analyzer.save_npz()`.

**Parameters** **fname** – The file's name.

#### **save** ()

Save state of analyzer, can be passed to `Analyzer.load()` to resemble the analyzer.

**Returns** The class name and the state.

#### **save\_npz** (*fname*)

Save state of analyzer, can be passed to `Analyzer.load_npz()` to resemble the analyzer.

**Parameters** **fname** – The file's name.

#### **class** `innvestigate.analyzer.base.TrainerMixin`

Mixin for analyzer that adapt to data.

This convenience interface exposes a Keras like training routing to the user.

**fit** (*X=None*, *batch\_size=32*, *\*\*kwargs*)

Takes the same parameters as Keras's `model.fit()` function.

**fit\_generator** (\**args*, *\*\*kwargs*)

Takes the same parameters as Keras's `model.fit_generator()` function.

**class** `innvestigate.analyzer.base.OneEpochTrainerMixin`

Exposes the same interface and functionality as `TrainerMixin` except that the training is limited to one epoch.

**fit** (\**args*, *\*\*kwargs*)

Same interface as `fit()` of `TrainerMixin` except that the parameter epoch is fixed to 1.

**fit\_generator** (\**args*, *\*\*kwargs*)

Same interface as `fit_generator()` of `TrainerMixin` except that the parameter epoch is fixed to 1.

**class** `innvestigate.analyzer.base.AnalyzerNetworkBase` (*model*, *neuron\_selection\_mode='max\_activation'*, *allow\_lambda\_layers=False*, *\*\*kwargs*)

Convenience interface for analyzers.

This class provides helpful functionality to create analyzer's. Basically it:

- takes the input model and adds a layer that selects the desired output neuron to analyze.
- passes the new model to `_create_analysis()` which should return the analysis as Keras tensors.
- compiles the function and serves the output to `analyze()` calls.
- allows `_create_analysis()` to return tensors that are intercept for debugging purposes.

#### Parameters

- **neuron\_selection\_mode** – How to select the neuron to analyze. Possible values are ‘max\_activation’, ‘index’ for the neuron (expects indices at `analyze()` calls), ‘all’ take all neurons.
- **allow\_lambda\_layers** – Allow the model to contain lambda layers.

**analyze** (*X*, *neuron\_selection=None*)

Same interface as Analyzer besides

**Parameters** **neuron\_selection** – If `neuron_selection_mode` is ‘index’ this should be an integer with the index for the chosen neuron.

**create\_analyzer\_model()**

Creates the `analyze` functionality. If not called beforehand it will be called by `analyze()`.

**class** `innvestigate.analyzer.base.ReverseAnalyzerBase` (*model*, *reverse\_verbose=False*, *reverse\_clip\_values=False*, *reverse\_project\_bottleneck\_layers=False*, *reverse\_check\_min\_max\_values=False*, *reverse\_check\_finite=False*, *reverse\_keep\_tensors=False*, *reverse\_reapply\_on\_copied\_layers=False*, *\*\*kwargs*)

Convenience class for analyzers that revert the model's structure.

This class contains many helper functions around the graph reverse function `innvestigate.utils.keras.graph.reverse_model()`.

The deriving classes should specify how the graph should be reverted by implementing the following functions:

- `_reverse_mapping(layer)` () given a layer this function returns a reverse mapping for the layer as specified in `innvestigate.utils.keras.graph.reverse_model()` or None.

This function can be implemented, but it is encouraged to implement a default mapping and add additional changes with the function `_add_conditional_reverse_mapping()` (see below).

The default behavior is finding a conditional mapping (see below), if none is found, `_default_reverse_mapping()` is applied.

- `_default_reverse_mapping()` defines the default reverse mapping.
- `_head_mapping()` defines how the outputs of the model should be instantiated before they are passed to the reversed network.

Furthermore other parameters of the function `innvestigate.utils.keras.graph.reverse_model()` can be changed by setting the according parameters of the init function:

#### Parameters

- `reverse_verbose` – Print information on the reverse process.
- `reverse_clip_values` – Clip the values that are passed along the reverted network. Expects tuple (min, max).
- `reverse_project_bottleneck_layers` – Project the value range of bottleneck tensors in the reverse network into another range.
- `reverse_check_min_max_values` – Print the min/max values observed in each tensor along the reverse network whenever `analyze()` is called.
- `reverse_check_finite` – Check if values passed along the reverse network are finite.
- `reverse_keep_tensors` – Keeps the tensors created in the backward pass and stores them in the attribute `_reversed_tensors`.
- `reverse_reapply_on_copied_layers` – See `innvestigate.utils.keras.graph.reverse_model()`.

## 1.2.2 Gradient methods

```
class innvestigate.analyzer.gradient_based.BaselineGradient(model, postprocess=None, **kwargs)
```

Gradient analyzer based on build-in gradient.

Returns as analysis the function value with respect to the input. The gradient is computed via the build in function. Is mainly used for debugging purposes.

**Parameters** `model` – A Keras model.

```
class innvestigate.analyzer.gradient_based.Gradient(model, postprocess=None, **kwargs)
```

Gradient analyzer.

Returns as analysis the function value with respect to the input. The gradient is computed via the library's network reverting.

**Parameters** `model` – A Keras model.

```
class innvestigate.analyzer.gradient_based.InputTimesGradient(model,
                                                               **kwargs)
Input*Gradient analyzer.

Parameters model – A Keras model.

class innvestigate.analyzer.gradient_based.Deconvnet(model, **kwargs)
Deconvnet analyzer.

Applies the “deconvnet” algorithm to analyze the model.

Parameters model – A Keras model.

class innvestigate.analyzer.gradient_based.GuidedBackprop(model, **kwargs)
Guided backprop analyzer.

Applies the “guided backprop” algorithm to analyze the model.

Parameters model – A Keras model.

class innvestigate.analyzer.gradient_based.IntegratedGradients(model,
                                                               steps=64,
                                                               **kwargs)
Integrated gradient analyzer.

Applies the “integrated gradient” algorithm to analyze the model.

Parameters
• model – A Keras model.
• steps – Number of steps to use average along integration path.

class innvestigate.analyzer.gradient_based.SmoothGrad(model, augment_by_n=64,
                                                       **kwargs)
Smooth grad analyzer.

Applies the “smooth grad” algorithm to analyze the model.

Parameters
• model – A Keras model.
• augment_by_n – Number of distortions to average for smoothing.
```

### 1.2.3 Layer-wise relevance propagation

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.BaselineLRPZ(model,
                                                                           **kwargs)
LRPZ analyzer - for testing purpose only.

Applies the “LRP-Z” algorithm to analyze the model. Based on the gradient times the input formula. This formula holds only for ReLU/MaxPooling networks, for which LRP-Z collapses into the stated formula.

Parameters model – A Keras model.

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRP(model,
                                                                    *args,
                                                                    **kwargs)
Base class for LRP-based model analyzers

Parameters
• model – A Keras model.
```

- **rule** – A rule can be a string or a Rule object, lists thereof or a list of conditions [(Condition, Rule), ...] gradient.

- **input\_layer\_rule** – either a Rule object, atuple of (low, high) the min/max pixel values of the inputs

**create\_rule\_mapping**(layer, reverse\_state)

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPZ(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-Z rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPZIgnoreBias(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-Z-ignore-bias rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPEpsilon(model,  
                           epsilon=1e-  
                           07,  
                           bias=True,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-Epsilon rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPEpsilonIgnoreBias(model,  
                           epsilon=1e-  
                           07,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-Epsilon-ignore-bias rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPWSquare(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the DeepTaylor W\*\*2 rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPFlat(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-Flat rule

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPAlphaBeta(model,  
                           al-  
                           pha=None,  
                           beta=None,  
                           bias=True,  
                           *args,  
                           **kwargs)
```

Base class for LRP AlphaBeta

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPAlpha2Beta1(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-alpha-beta rule with a=2,b=1

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPAlpha2Beta1IgnoreBias(model,  
                           *args,  
                           **kwargs)
```

LRP-analyzer that uses the LRP-alpha-beta-ignbias rule with a=2,b=1

```
class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPAlpha1Beta0(model,
*args,
**kwargs)
    LRP-analyzer that uses the LRP-alpha-beta rule with a=1,b=0

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPAlpha1Beta0IgnoreBias(model,
*args,
**kwargs)
    LRP-analyzer that uses the LRP-alpha-beta-ignbias rule with a=1,b=0

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPZPlus(model,
*args,
**kwargs)
    LRP-analyzer that uses the LRP-alpha-beta rule with a=1,b=0

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPZPlusFast(model,
*args,
**kwargs)
    The ZPlus rule is a special case of the AlphaBetaRule for alpha=1, beta=0 and assumes inputs x >= 0.

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPSequentialPresetA(model,
ep-
silon=0.1
*args,
**kwargs)
    Special LRP-configuration for ConvNets

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPSequentialPresetB(model,
ep-
silon=0.1
*args,
**kwargs)
    Special LRP-configuration for ConvNets

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPSequentialPresetAFlat(model,
*args,
**kwargs)
    Special LRP-configuration for ConvNets

class innvestigate.analyzer.relevance_based.relevance_analyzer.LRPSequentialPresetBFlat(model,
*args,
**kwargs)
    Special LRP-configuration for ConvNets
```

## 1.2.4 Pattern methods

```
class innvestigate.analyzer.pattern_based.PatternNet(model, patterns=None, pattern_type=None, **kwargs)
```

PatternNet analyzer.

Applies the “PatternNet” algorithm to analyze the model’s predictions.

### Parameters

- **model** – A Keras model.
- **patterns** – Pattern computed by innvestigate.tools.PatternComputer. If None fit() needs to be called.
- **allow\_lambda\_layers** – Approximate lambda layers with the gradient.

- **reverse\_project\_bottleneck\_layers** – Project the analysis vector into range [-1, +1]. (default: True)

```
class innvestigate.analyzer.pattern_based.PatternAttribution(model,          patterns=None,      pattern_type=None,      **kwargs)
```

PatternAttribution analyzer.

Applies the “PatternNet” algorithm to analyze the model’s predictions.

#### Parameters

- **model** – A Keras model.
- **patterns** – Pattern computed by `innvestigate.tools.PatternComputer`. If `None` `fit()` needs to be called.
- **allow\_lambda\_layers** – Approximate lambda layers with the gradient.
- **reverse\_project\_bottleneck\_layers** – Project the analysis vector into range [-1, +1]. (default: True)

## 1.2.5 Deep Taylor

```
class innvestigate.analyzer.deeptaylor.DeepTaylor(model, *args, **kwargs)
```

DeepTaylor for ReLU-networks with unbounded input

This class implements the DeepTaylor algorithm for neural networks with ReLU activation and unbounded input ranges.

#### Parameters **model** – A Keras model.

```
class innvestigate.analyzer.deeptaylor.BoundedDeepTaylor(model,           low=None,          high=None,      **kwargs)
```

DeepTaylor for ReLU-networks with bounded input

This class implements the DeepTaylor algorithm for neural networks with ReLU activation and bounded input ranges.

#### Parameters

- **model** – A Keras model.
- **low** – Lowest value of the input range. See Z\_B rule.
- **high** – Highest value of the input range. See Z\_B rule.

## 1.2.6 DeepLIFT

```
class innvestigate.analyzer.deeplift.DeepLIFTWrapper(model, **kwargs)
```

Wrapper around DeepLIFT package

This class wraps the DeepLIFT package. For further explanation of the parameters check out: <https://github.com/kundajelab/deeplift>

#### Parameters

- **model** – A Keras model.
- **nonlinear\_mode** – The nonlinear mode parameter.
- **reference\_inputs** – The reference input used for DeepLIFT.

- **verbose** – Verbosity of the DeepLIFT package.

**Note** Requires the deeplift package.

**analyze** (*X*, *neuron\_selection=None*)  
Same interface as Analyzer besides

**Parameters** **neuron\_selection** – If *neuron\_selection\_mode* is ‘index’ this should be an integer with the index for the chosen neuron.

## 1.2.7 Misc

**class** innvestigate.analyzer.misc.**Random** (*model*, *stddev=1*, *\*\*kwargs*)  
Returns noise.

Returns the Gaussian noise as analysis.

**Parameters**

- **model** – A Keras model.
- **stddev** – The standard deviation of the noise.

**class** innvestigate.analyzer.misc.**Input** (*model*, *neuron\_selection\_mode='max\_activation'*,  
*allow\_lambda\_layers=False*, *\*\*kwargs*)

Returns the input.

Returns the input as analysis.

**Parameters** **model** – A Keras model.

## 1.3 innvestigate.applications

### 1.3.1 Imagenet

Example applications for image classification.

Each function returns a pretrained ImageNet model. The models are based on keras.applications models and contain additionally pretrained patterns.

The returned dictionary contains the following keys: *model*, *in*, *sm\_out*, *out*, *image\_shape*, *color\_coding*, *preprocess\_f*, *patterns*.

Function parameters:

**param** **load\_weights** Download or access cached weights.

**param** **load\_patterns** Download or access cached patterns.

```
innvestigate.applications.imagenet.vgg16 (load_weights=False, load_patterns=False)
innvestigate.applications.imagenet.vgg19 (load_weights=False, load_patterns=False)
innvestigate.applications.imagenet.resnet50 (load_weights=False, load_patterns=False)
innvestigate.applications.imagenet.inception_v3 (load_weights=False,
                                              load_patterns=False)
innvestigate.applications.imagenet.inception_resnet_v2 (load_weights=False,
                                                       load_patterns=False)
```

```
innvestigate.applications.imagenet.densenet121 (load_weights=False,  
                                load_patterns=False)  
innvestigate.applications.imagenet.densenet169 (load_weights=False,  
                                load_patterns=False)  
innvestigate.applications.imagenet.densenet201 (load_weights=False,  
                                load_patterns=False)  
innvestigate.applications.imagenet.nasnet_large (load_weights=False,  
                                load_patterns=False)  
innvestigate.applications.imagenet.nasnet_mobile (load_weights=False,  
                                load_patterns=False)
```

## 1.4 innvestigate.tools

### 1.4.1 Pattern Computation

```
class innvestigate.tools.pattern.PatternComputer (model, pattern_type='linear', com-  
                                         pute_layers_in_parallel=True,  
                                         gpus=None)
```

Pattern computer.

Computes a pattern for each layer with a kernel of a given model.

#### Parameters

- **model** – A Keras model.
- **pattern\_type** – A string or a tuple of strings. Valid types are ‘linear’, ‘relu’, ‘relu.positive’, ‘relu.negative’.
- **compute\_layers\_in\_parallel** – Not supported yet. Compute all patterns at once. Otherwise computer layer after layer.
- **gpus** – Not supported yet. Gpus to use.

```
compute (X, batch_size=32, verbose=0)
```

Compute and return the patterns for the model and the data *X*.

#### Parameters

- **X** – Data to compute patterns.
- **batch\_size** – Batch size to use.
- **verbose** – As for keras model.fit.

```
compute_generator (generator, **kwargs)
```

Compute and return the patterns for the model and the data *X*.

#### Parameters

- **generator** – Data to compute patterns.
- **kwargs** – Same as for keras model.fit\_generator.

## 1.4.2 Perturbation Analysis

```
class innvestigate.tools.perturbate.Perturbation(perturbation_function,
                                                num_perturbed_regions=0,
                                                region_shape=(9, 9),      re-
                                                duce_function=<function
                                                mean>,                  aggrega-
                                                tion_function=<function    mean>,
                                                pad_mode='reflect', in_place=False,
                                                value_range=None)
```

Perturbation of pixels based on analysis result.

### Parameters

- **perturbation\_function** (*function or callable or str*) – Defines the function with which the samples are perturbated. Can be a function or a string that defines a predefined perturbation function.
- **num\_perturbed\_regions** (*int*) – Number of regions to be perturbed.
- **reduce\_function** (*function or callable*) – Function to reduce the analysis result to one channel, e.g. mean or max function.
- **aggregation\_function** (*function or callable*) – Function to aggregate the analysis over subregions.
- **pad\_mode** (*str or function or callable*) – How to pad if the image cannot be subdivided into an integer number of regions. As in numpy.pad.
- **in\_place** (*bool*) – If true, the perturbations are performed in place, i.e. the input samples are modified.
- **value\_range** (*tuple*) – Minimal and maximal value after perturbation as a tuple: (min\_val, max\_val). The input is clipped to this range

```
aggregate_regions(analysis)

static compute_perturbation_mask(ranks, num_perturbated_regions)
static compute_region_ordering(aggregated_regions)
expand_regions_to_pixels(regions)

pad(analysis)

perturbate_on_batch(x, analysis)
```

### Parameters

- **x** (*numpy.ndarray*) – Batch of images.
- **analysis** (*numpy.ndarray*) – Analysis of this batch.

**Returns** Batch of perturbated images

**Return type** numpy.ndarray

```
perturbate_regions(x, perturbation_mask_regions)
reshape_region_pixels(region_pixels, target_shape)
reshape_to_regions(analysis)
```

```
class innvestigate.tools.perturbate.PerturbationAnalysis(analyzer, model,
                                                       generator, pertur-
                                                       bation, steps=1, re-
                                                       gions_per_step=1, re-
                                                       compute_analysis=False,
                                                       verbose=False)
```

Performs the perturbation analysis.

#### Parameters

- **analyzer** (`innvestigate.analyzer.base.AnalyzerBase`) – Analyzer.
- **model** (`keras.engine.training.Model`) – Trained Keras model.
- **generator** (`innvestigate.utils.BatchSequence`) – Data generator.
- **perturbation** (`innvestigate.tools.Perturbation`) – Instance of Perturbation class that performs the perturbation.
- **steps** (`int`) – Number of perturbation steps.
- **regions\_per\_step** (`float`) – Number of regions that are perturbed per step.
- **recompute\_analysis** (`bool`) – If true, the analysis is recomputed after each perturbation step.
- **verbose** – If true, print some useful information, e.g. timing, progress etc.

```
compute_on_batch(x, analysis=None, return_analysis=False)
```

Computes the analysis and perturbs the input batch accordingly.

#### Parameters

- **x** (`numpy.ndarray`) – Samples.
- **analysis** – Analysis of x. If None, it is recomputed.

```
compute_perturbation_analysis()
```

```
evaluate_generator(generator, steps=None, max_queue_size=10, workers=1,
                   use_multiprocessing=False)
```

Evaluates the model on a data generator.

The generator should return the same kind of data as accepted by `test_on_batch`. For documentation, refer to `keras.engine.training.evaluate_generator` (<https://keras.io/models/model/>)

```
evaluate_on_batch(x, y, analysis=None, sample_weight=None)
```

Perturbs the input batch and scores the model on the perturbed batch.

#### Parameters

- **x** (`numpy.ndarray`) – Samples.
- **y** (`numpy.ndarray`) – Labels.
- **analysis** (`numpy.ndarray`) – Analysis of x.
- **sample\_weight** (`None`) – Sample weights.

**Returns** List of test scores.

**Return type** list

## 1.5 innvestigate.utils

```
innvestigate.utils.model_wo_softmax(*args, **kwargs)
```

```
innvestigate.utils.to_list(l)
```

If not list, wraps parameter into a list.

```
class innvestigate.utils.BatchSequence(Xs, batch_size=32)
```

Batch sequence generator.

Take a (list of) input tensors and a batch size and creates a generators that creates a sequence of batches.

### Parameters

- **Xs** – One or a list of tensors. First axis needs to have same length.
- **batch\_size** – Batch size. Default 32.

```
class innvestigate.utils.TargetAugmentedSequence(sequence, augment_f)
```

Augments a sequence with a target on the fly.

Takes a sequence/generator and a function that creates on the fly for each batch a target. The generator takes a batch from that sequence, computes the target and returns both.

### Parameters

- **sequence** – A sequence or generator.
- **augment\_f** – Takes a batch and returns a target.

```
innvestigate.utils.preprocess_images(images, color_coding=None)
```

Image preprocessing

Takes a batch of images and: \* Adjust the color axis to the Keras format. \* Fixes the color coding.

### Parameters

- **images** – Batch of images with 4 axes.
- **color\_coding** – Determines the color coding. Can be None, ‘RGBtoBGR’ or ‘BGR-toRGB’.

**Returns** The preprocessed batch.

```
innvestigate.utils.postprocess_images(images, color_coding=None, channels_first=None)
```

Image postprocessing

Takes a batch of images and reverts the preprocessing.

### Parameters

- **images** – A batch of images with 4 axes.
- **color\_coding** – The initial color coding, see [preprocess\\_images\(\)](#).
- **channels\_first** – The output channel format.

**Returns** The postprocessed images.

## 1.5.1 Visualizations

```
innvestigate.utils.visualizations.project(X, output_range=(0, 1), absmax=None, input_is_positive_only=False)
```

Projects a tensor into a value range.

Projects the tensor values into the specified range.

#### Parameters

- **x** – A tensor.
- **output\_range** – The output value range.
- **absmax** – A tensor specifying the absmax used for normalizing. Default the absmax along the first axis.
- **input\_is\_positive\_only** – Is the input value range only positive.

**Returns** The tensor with the values project into output range.

```
innvestigate.utils.visualizations.heatmap(X, cmap_type='seismic', reduce_op='sum',
                                         reduce_axis=-1, alpha_cmap=False,
                                         **kwargs)
```

Creates a heatmap/color map.

Create a heatmap or colormap out of the input tensor.

#### Parameters

- **x** – A image tensor with 4 axes.
- **cmap\_type** – The color map to use. Default ‘seismic’.
- **reduce\_op** – Operation to reduce the color axis. Either ‘sum’ or ‘absmax’.
- **reduce\_axis** – Axis to reduce.
- **alpha\_cmap** – Should the alpha component of the cmap be included.
- **kwargs** – Arguments passed on to *project ()*

**Returns** The tensor as color-map.

```
innvestigate.utils.visualizations.graymap(X, **kwargs)
```

Same as *heatmap ()* but uses a gray colormap.

```
innvestigate.utils.visualizations.gamma(X, gamma=0.5, minamp=0, maxamp=None)
```

Apply gamma correction to an input array X while maintaining the relative order of entries, also for negative vs positive values in X. the fxn firstly determines the max amplitude in both positive and negative direction and then applies gamma scaling to the positive and negative values of the array separately, according to the common amplitude.

#### Parameters

- **gamma** – the gamma parameter for gamma scaling
- **minamp** – the smallest absolute value to consider. if not given assumed to be zero (neutral value for relevance, min value for saliency, ...). values above and below minamp are treated separately.
- **maxamp** – the largest absolute value to consider relative to the neutral value minamp if not given determined from the given data.

```
innvestigate.utils.visualizations.clip_quantile(X, quantile=1)
```

Clip the values of X into the given quantile.

## 1.5.2 Keras-Utils

```
innvestigate.utils.keras.graph.copy_layer(layer, keep_bias=True, name_template=None,  
weights=None, reuse_symbolic_tensors=True,  
**kwargs)
```

Copy a Keras layer

Copies a Keras layer.

### Parameters

- **layer** – A layer that should be copied.
- **keep\_bias** – Keep a potential bias.
- **weights** – Weights to set in the new layer. Options: np tensors, symbolic tensors, or None, in which case the weights from old\_layer are used.
- **reuse\_symbolic\_tensors** – If the weights of the old\_layer are used copy the symbolic ones or copy the Numpy weights.

**Returns** The new layer instance.

```
innvestigate.utils.keras.graph.model_wo_softmax(model)
```

Creates a new model w/o the final softmax activation.

```
innvestigate.utils.keras.graph.get_model_execution_graph(model,
```

*keep\_input\_layers=False*

Returns a dictionary representing the execution graph. Each key is the node's id as it is used by the reverse\_model method.

Each associated value contains a dictionary with the following items:

- nid: the node id.
- layer: the layer creating this node.
- Xs: the input tensors (only valid if not in a nested container).
- Ys: the output tensors (only valid if not in a nested container).
- Xs\_nids: the ids of the nodes creating the Xs.
- Ys\_nids: the ids of nodes using the according output tensor.
- Xs\_layers: the layer that created the according input tensor.
- Ys\_layers: the layers using the according output tensor.

### Parameters

- **model** – A kera model.
- **keep\_input\_layers** – Keep input layers.

```
innvestigate.utils.keras.graph.reverse_model (model,           reverse_mappings,      de-
                                              fault_reverse_mapping=None,
                                              head_mapping=None,
                                              stop_mapping_at_tensors=[],
                                              verbose=False,          re-
                                              turn_all_reversed_tensors=False,
                                              clip_all_reversed_tensors=False,
                                              project_bottleneck_tensors=False,
                                              execution_trace=None,      reap-
                                              ply_on_copied_layers=False)
```

Reverses a Keras model based on the given reverse functions. It returns the reverted tensors for the according model inputs.

### Parameters

- **model** – A Keras model.
- **reverse\_mappings** – Either a callable that matches layers to mappings or a dictionary with layers as keys and mappings as values. Allowed as mapping forms are:
  - A function of form (A) f(Xs, Ys, reversed\_Ys, reverse\_state).
  - A function of form f(B) f(layer, reverse\_state) that returns a function of form (A).
  - A ReverseMappingBase subclass.
- **default\_reverse\_mapping** – A function that reverses layers for which no mapping was given by param “reverse\_mappings”.
- **head\_mapping** – Map output tensors to new values before passing them into the reverted network.
- **stop\_mapping\_at\_tensors** – Tensors at which to stop the mapping. Similar to stop\_gradient parameters for gradient computation.
- **verbose** – Print what’s going on.
- **return\_all\_reversed\_tensors** – Return all reverted tensors in addition to reverted model input tensors.
- **clip\_all\_reversed\_tensors** – Clip each reverted tensor. False or tuple with min/max value.
- **project\_bottleneck\_tensors** – Project bottleneck layers in the reverting process into a given value range. False, True or (a, b) for projection range.
- **reapply\_on\_copied\_layers** – When a model execution needs to linearized and copy layers before reapplying them. See `trace_model_execution()`.

# CHAPTER 2

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Python Module Index

---

i

innvestigate.analyzer.base, 3  
innvestigate.analyzer.deeplift, 10  
innvestigate.analyzer.deeptaylor, 10  
innvestigate.analyzer.gradient\_based, 6  
innvestigate.analyzer.misc, 11  
innvestigate.analyzer.pattern\_based, 9  
innvestigate.analyzer.relevance\_based.relevance\_analyzer,  
    7  
innvestigate.applications.imagenet, 11  
innvestigate.tools.pattern, 12  
innvestigate.tools.perturbate, 13  
innvestigate.utils, 15  
innvestigate.utils.keras.graph, 17  
innvestigate.utils.visualizations, 15



---

## Index

---

### A

aggregate\_regions() (innvestigate.tools.perturbate.Perturbation method), 13  
analyze() (innvestigate.analyzer.base.AnalyzerBase method), 4  
analyze() (innvestigate.analyzer.base.AnalyzerNetworkBase method), 5  
analyze() (innvestigate.analyzer.deeplift.DeepLIFTWrapper method), 11  
AnalyzerBase (class in innvestigate.analyzer.base), 3  
AnalyzerNetworkBase (class in innvestigate.analyzer.base), 5

### B

BaselineGradient (class in innvestigate.analyzer.gradient\_based), 6  
BaselineLRPZ (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 7  
BatchSequence (class in innvestigate.utils), 15  
BoundedDeepTaylor (class in innvestigate.analyzer.deeptaylor), 10

### C

clip\_quantile() (in module gate.utils.visualizations), 16  
compute() (innvestigate.tools.pattern.PatternComputer method), 12  
compute\_generator() (innvestigate.tools.pattern.PatternComputer method), 12  
compute\_on\_batch() (innvestigate.tools.perturbate.PerturbationAnalysis method), 14

compute\_perturbation\_analysis() (innvestigate.tools.perturbate.PerturbationAnalysis method), 14  
compute\_perturbation\_mask() (innvestigate.tools.perturbate.Perturbation method), 13  
compute\_region\_ordering() (innvestigate.tools.perturbate.Perturbation method), 13  
copy\_layer() (in module gate.utils.keras.graph), 17  
create\_analyzer() (in module innvestigate.analyzer), 3  
create\_analyzer\_model() (innvestigate.analyzer.base.AnalyzerNetworkBase method), 5  
create\_rule\_mapping() (innvestigate.analyzer.relevance\_based.relevance\_analyzer.LRP method), 8

### D

Deconvnet (class in innvestigate.analyzer.gradient\_based), 7  
DeepLIFTWrapper (class in innvestigate.deeplift), 10  
DeepTaylor (class in innvestigate.analyzer.deeptaylor), 10  
densenet121() (in module gate.applications.imagenet), 11  
densenet169() (in module gate.applications.imagenet), 12  
densenet201() (in module gate.applications.imagenet), 12

### E

evaluate\_generator() (innvestigate.tools.perturbate.PerturbationAnalysis method), 14  
evaluate\_on\_batch() (innvestigate.tools.perturbate.PerturbationAnalysis

method), 14  
expand\_regions\_to\_pixels() (innvestigate.tools.perturbate.Perturbation method), 13

**F**

fit() (innvestigate.analyzer.base.AnalyzerBase method), 4  
fit() (innvestigate.analyzer.base.OneEpochTrainerMixin method), 5  
fit() (innvestigate.analyzer.base.TrainerMixin method), 5  
fit\_generator() (innvestigate.analyzer.base.AnalyzerBase 4)  
fit\_generator() (innvestigate.analyzer.base.OneEpochTrainerMixin method), 5  
fit\_generator() (innvestigate.analyzer.base.TrainerMixin 5)

**G**

gamma() (in module innvestigate.utils.visualizations), 16  
get\_model\_execution\_graph() (in module innvestigate.utils.keras.graph), 17  
Gradient (class in innvestigate.analyzer.gradient\_based), 6  
graymap() (in module innvestigate.utils.visualizations), 16  
GuidedBackprop (class in innvestigate.analyzer.gradient\_based), 7

**H**

heatmap() (in module innvestigate.utils.visualizations), 16

**I**

inception\_resnet\_v2() (in module innvestigate.applications.imagenet), 11  
inception\_v3() (in module innvestigate.applications.imagenet), 11  
innvestigate.analyzer.base (module), 3  
innvestigate.analyzer.deeplift (module), 10  
innvestigate.analyzer.deeptaylor (module), 10  
innvestigate.analyzer.gradient\_based (module), 6  
innvestigate.analyzer.misc (module), 11  
innvestigate.analyzer.pattern\_based (module), 9

innvestigate.analyzer.relevance\_based.relevance\_analyzer (module), 7  
innvestigate.applications.imagenet (module), 11  
innvestigate.tools.pattern (module), 12  
innvestigate.tools.perturbate (module), 13  
innvestigate.utils (module), 15  
innvestigate.utils.keras.graph (module), 17  
innvestigate.utils.visualizations (module), 15  
Input (class in innvestigate.analyzer.misc), 11  
InputTimesGradient (class in innvestigate.analyzer.gradient\_based), 6  
IntegratedGradients (class in innvestigate.analyzer.gradient\_based), 7

**L**

load() (innvestigate.analyzer.base.AnalyzerBase static method), 4  
load\_npz() (innvestigate.analyzer.base.AnalyzerBase static method), 4  
LRP (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 7  
LRPAlphaBeta0 (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPAlphaBeta0IgnoreBias (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 9  
LRPAlpha2Beta1 (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPAlpha2Beta1IgnoreBias (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPAlphaBeta (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPEpsilon (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPEpsilonIgnoreBias (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPFlat (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 8  
LRPSequentialPresetA (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 9  
LRPSequentialPresetAFlat (class in innvestigate.analyzer.relevance\_based.relevance\_analyzer), 9

|   |   |                  |
|---|---|------------------|
| 9   | perturbate_regions ()   | (investigate), 9 |
| LRPSequentialPresetB (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 9     | gate.tools.perturbate.Perturbation                                    | (method), 13     |
| 9   | Perturbation (class in investigate.tools.perturbate), 13              |                  |
| LRPSequentialPresetBFlat (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 9 | PerturbationAnalysis (class in investigate.gate.tools.perturbate), 13 |                  |
| 9   | postprocess_images () (in module investigate.gate.utils), 15          |                  |
| LRPWSquare (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 8               | preprocess_images () (in module investigate.gate.utils), 15           |                  |
| 8   | project () (in module investigate.gate.utils.visualizations), 15      |                  |
| 8   | R   |                  |
| LRPZIgnoreBias (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 8           | Random (class in investigate.analyzer.misc), 11                       |                  |
| LRPZPlus (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 9                 | reshape_region_pixels ()  | (method), 13     |
| 9   | gate.tools.perturbate.Perturbation                                    |                  |
| LRPZPlusFast (class in investigate.gate.analyzer.relevance_based.relevance_analyzer), 9             | reshape_to_regions ()   | (method), 13     |
| 9   | gate.tools.perturbate.Perturbation                                    |                  |
| M   | resnet50 () (in module investigate.gate.applications.imagenet), 11    |                  |
| model_wo_softmax () (in module investigate.utils), 15   | reverse_model () (in module investigate.gate.utils.keras.graph), 17   |                  |
| model_wo_softmax () (in module investigate.gate.utils.keras.graph), 17                              | ReverseAnalyzerBase (class in investigate.analyzer.base), 5           |                  |
| N   | S   |                  |
| nasnet_large () (in module investigate.gate.applications.imagenet), 12                              | save () (investigate.analyzer.base.AnalyzerBase method), 4            |                  |
| nasnet_mobile () (in module investigate.gate.applications.imagenet), 12                             | save_npz () (investigate.analyzer.base.AnalyzerBase method), 4        |                  |
| NotAnalyzeableModelError, 3   | SmoothGrad (class in investigate.analyzer.gradient_based), 7          |                  |
| O   | T   |                  |
| OneEpochTrainerMixin (class in investigate.gate.analyzer.base), 5                                   | TargetAugmentedSequence (class in investigate.gate.utils), 15         |                  |
| P   | to_list () (in module investigate.utils), 15                          |                  |
| pad () (investigate.tools.perturbate.Perturbation method), 13                                       | TrainerMixin (class in investigate.analyzer.base), 4                  |                  |
| PatternAttribution (class in investigate.gate.analyzer.pattern_based), 10                           | V   |                  |
| PatternComputer (class in investigate.gate.tools.pattern), 12                                       | vgg16 () (in module investigate.gate.applications.imagenet), 11       |                  |
| PatternNet (class in investigate.gate.analyzer.pattern_based), 9                                    | vgg19 () (in module investigate.gate.applications.imagenet), 11       |                  |
| perturbate_on_batch ()  | (investigate.tools.perturbate.Perturbation method), 13                |                  |