
Foolbox Documentation

Release 1.9.0

Jonas Rauber & Wieland Brendel

Apr 04, 2019

1	Robust Vision Benchmark	3
1.1	Installation	3
1.2	Tutorial	4
1.3	Examples	5
1.4	Advanced	9
1.5	Model Zoo	10
1.6	Development	10
1.7	FAQ	11
1.8	foolbox.models	11
1.9	foolbox.criteria	28
1.10	foolbox.zoo	33
1.11	foolbox.distances	35
1.12	foolbox.attacks	36
1.13	foolbox.adversarial	59
1.14	foolbox.utils	62
2	Indices and tables	65
	Bibliography	67
	Python Module Index	69

Foolbox is a Python toolbox to create adversarial examples that fool neural networks.

It comes with support for many frameworks to build models including

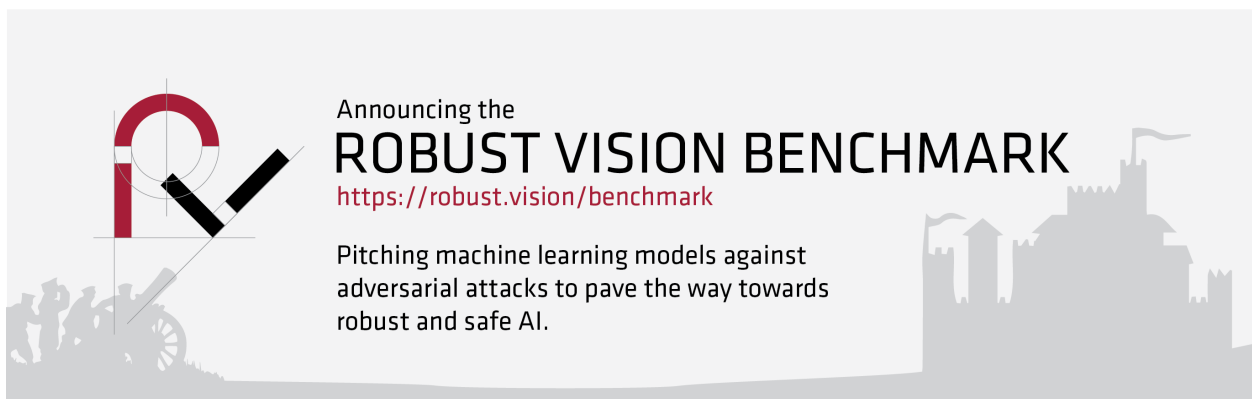
- TensorFlow
- PyTorch
- Theano
- Keras
- Lasagne
- MXNet

and it is easy to extend to other frameworks.

In addition, it comes with a **large collection of adversarial attacks**, both gradient-based attacks as well as black-box attacks. See [*foolbox.attacks*](#) for details.

The source code and a [minimal working example](#) can be found on [GitHub](#).

Robust Vision Benchmark



You might want to have a look at our recently announced [Robust Vision Benchmark](https://robust.vision/benchmark), a benchmark for adversarial attacks and the robustness of machine learning models.

1.1 Installation

Foolbox is a Python package to create adversarial examples. We test using Python 2.7, 3.5 and 3.6, but other versions of Python might work as well. **We recommend using Python 3!**

1.1.1 Stable release

You can install the latest stable release of Foolbox from PyPI using *pip*:

```
pip install foolbox
```

Make sure that *pip* installs packages for Python 3, otherwise you might need to use *pip3* instead of *pip*.

1.1.2 Development version

Alternatively, you can install the latest development version of Foolbox from GitHub. We try to keep the master branch stable, so this version should usually work fine. Feel free to open an issue on GitHub if you encounter any problems.

```
pip install https://github.com/bethgelab/foolbox/archive/master.zip
```

1.1.3 Contributing to Foolbox

If you would like to contribute the development of Foolbox, install it in editable mode:

```
git clone https://github.com/bethgelab/foolbox.git
cd foolbox
pip install --editable .
```

To contribute your changes, you will need to fork the Foolbox repository on GitHub. You can then add it as a remote:

```
git remote rename origin upstream
git remote add origin https://github.com/<your-github-name>/foolbox.git
```

You can now commit your changes, push them to your fork and create a pull-request to contribute them to Foolbox.

1.2 Tutorial

This tutorial will show you how an adversarial attack can be used to find adversarial examples for a model.

1.2.1 Creating a model

For the tutorial, we will target *VGG19* implemented in *TensorFlow*, but it is straight forward to apply the same to other models or other frameworks such as *Theano* or *PyTorch*.

```
import tensorflow as tf

images = tf.placeholder(tf.float32, (None, 224, 224, 3))
preprocessed = vgg_preprocessing(images)
logits = vgg19(preprocessed)
```

To turn a model represented as a standard TensorFlow graph into a model that can be attacked by the Adversarial Toolbox, all we have to do is to create a new *TensorFlowModel* instance:

```
from foolbox.models import TensorFlowModel

model = TensorFlowModel(images, logits, bounds=(0, 255))
```

1.2.2 Specifying the criterion

To run an adversarial attack, we need to specify the type of adversarial we are looking for. This can be done using the *Criterion* class.


```

from foolbox.criteria import TargetClassProbability

target_class = 22
criterion = TargetClassProbability(target_class, p=0.99)

```

1.2.3 Running the attack

Finally, we can create and apply the attack:

```

from foolbox.attacks import LBFGSAttack

attack = LBFGSAttack(model, criterion)

image = np.asarray(Image.open('example.jpg'))
label = np.argmax(model.predictions(image))

adversarial = attack(image, label=label)

```

1.2.4 Visualizing the adversarial examples

To plot the adversarial example we can use *matplotlib*:

```

import matplotlib.pyplot as plt

plt.subplot(1, 3, 1)
plt.imshow(image)

plt.subplot(1, 3, 2)
plt.imshow(adversarial)

plt.subplot(1, 3, 3)
plt.imshow(adversarial - image)

```

1.3 Examples

Here you can find a collection of examples how Foolbox models can be created using different deep learning frameworks and some full-blown attack examples at the end.

1.3.1 Creating a model

Keras: ResNet50

```

import keras
import numpy as np
import foolbox

keras.backend.set_learning_phase(0)
kmodel = keras.applications.resnet50.ResNet50(weights='imagenet')
preprocessing = (np.array([104, 116, 123]), 1)

```

(continues on next page)

(continued from previous page)

```

model = foolbox.models.KerasModel(kmodel, bounds=(0, 255),
↳preprocessing=preprocessing)

image, label = foolbox.utils.imagenet_example()
# ::-1 reverses the color channels, because Keras ResNet50 expects BGR instead of RGB
print(np.argmax(model.predictions(image[:, :, ::-1])), label)

```

PyTorch: ResNet18

You might be interested in checking out the full PyTorch example at the end of this document.

```

import torchvision.models as models
import numpy as np
import foolbox

# instantiate the model
resnet18 = models.resnet18(pretrained=True).cuda().eval() # for CPU, remove cuda()
mean = np.array([0.485, 0.456, 0.406]).reshape((3, 1, 1))
std = np.array([0.229, 0.224, 0.225]).reshape((3, 1, 1))
model = foolbox.models.PyTorchModel(resnet18, bounds=(0, 1), num_classes=1000,
↳preprocessing=(mean, std))

image, label = foolbox.utils.imagenet_example(data_format='channels_first')
image = image / 255
print(np.argmax(model.predictions(image)), label)

```

TensorFlow: VGG19

First, create the model in TensorFlow.

```

import tensorflow as tf
from tensorflow.contrib.slim.nets import vgg
import numpy as np
import foolbox

images = tf.placeholder(tf.float32, shape=(None, 224, 224, 3))
preprocessed = images - [123.68, 116.78, 103.94]
logits, _ = vgg.vgg_19(preprocessed, is_training=False)
restorer = tf.train.Saver(tf.trainable_variables())

image, _ = foolbox.utils.imagenet_example()

```

Then transform it into a Foolbox model using one of these four options:

Option 1

This option is recommended if you want to keep the code as short as possible. It makes use of the TensorFlow session created by Foolbox internally if no default session is set.

```

with foolbox.models.TensorFlowModel(images, logits, (0, 255)) as model:
    restorer.restore(model.session, '/path/to/vgg_19.ckpt')
    print(np.argmax(model.predictions(image)))

```

Option 2

This option is recommended if you want to create the TensorFlow session yourself.

```
with tf.Session() as session:
    restorer.restore(session, '/path/to/vgg_19.ckpt')
    model = foolbox.models.TensorFlowModel(images, logits, (0, 255))
    print(np.argmax(model.predictions(image)))
```

Option 3

This option is recommended if you want to avoid nesting context managers, e.g. during interactive development.

```
session = tf.InteractiveSession()
restorer.restore(session, '/path/to/vgg_19.ckpt')
model = foolbox.models.TensorFlowModel(images, logits, (0, 255))
print(np.argmax(model.predictions(image)))
session.close()
```

Option 4

This is possible, but usually one of the other options should be preferred.

```
session = tf.Session()
with session.as_default():
    restorer.restore(session, '/path/to/vgg_19.ckpt')
    model = foolbox.models.TensorFlowModel(images, logits, (0, 255))
    print(np.argmax(model.predictions(image)))
session.close()
```

1.3.2 Applying an attack

Once you created a Foolbox model (see the previous section), you can apply an attack.

FGSM (GradientSignAttack)

```
# create a model (see previous section)
fmodel = ...

# get source image and label
image, label = foolbox.utils.imagenet_example()

# apply attack on source image
attack = foolbox.attacks.FGSM(fmodel)
adversarial = attack(image[:, :, ::-1], label)
```

1.3.3 Creating an untargeted adversarial for a PyTorch model

```
import foolbox
import torch
import torchvision.models as models
import numpy as np

# instantiate the model
resnet18 = models.resnet18(pretrained=True).eval()
if torch.cuda.is_available():
    resnet18 = resnet18.cuda()
mean = np.array([0.485, 0.456, 0.406]).reshape((3, 1, 1))
std = np.array([0.229, 0.224, 0.225]).reshape((3, 1, 1))
fmodel = foolbox.models.PyTorchModel(
    resnet18, bounds=(0, 1), num_classes=1000, preprocessing=(mean, std))

# get source image and label
image, label = foolbox.utils.imagenet_example(data_format='channels_first')
image = image / 255. # because our model expects values in [0, 1]

print('label', label)
print('predicted class', np.argmax(fmodel.predictions(image)))

# apply attack on source image
attack = foolbox.attacks.FGSM(fmodel)
adversarial = attack(image, label)

print('adversarial class', np.argmax(fmodel.predictions(adversarial)))
```

outputs

```
label 282
predicted class 282
adversarial class 281
```

To plot image and adversarial, don't forget to move the channel axis to the end before passing them to matplotlib's `imshow`, e.g. using `np.transpose(image, (1, 2, 0))`.

1.3.4 Creating a targeted adversarial for the Keras ResNet model

```
import foolbox
from foolbox.models import KerasModel
from foolbox.attacks import LBFGSAttack
from foolbox.criteria import TargetClassProbability
import numpy as np
import keras
from keras.applications.resnet50 import ResNet50
from keras.applications.resnet50 import preprocess_input
from keras.applications.resnet50 import decode_predictions

keras.backend.set_learning_phase(0)
kmodel = ResNet50(weights='imagenet')
preprocessing = (np.array([104, 116, 123]), 1)
fmodel = KerasModel(kmodel, bounds=(0, 255), preprocessing=preprocessing)

image, label = foolbox.utils.imagenet_example()
```

(continues on next page)

(continued from previous page)

```
# run the attack
attack = LBFGSAttack(model=fmodel, criterion=TargetClassProbability(781, p=.5))
adversarial = attack(image[:, :, :-1], label)

# show results
print(np.argmax(fmodel.predictions(adversarial)))
print(foolbox.utils.softmax(fmodel.predictions(adversarial))[781])
adversarial_rgb = adversarial[np.newaxis, :, :, :-1]
preds = kmodel.predict(preprocess_input(adversarial_rgb.copy()))
print("Top 5 predictions (adversarial: ", decode_predictions(preds, top=5))
```

outputs

```
781
0.832095
Top 5 predictions (adversarial: [(('n04149813', 'scoreboard', 0.83013469), (
↪ 'n03196217', 'digital_clock', 0.030192226), ('n04152593', 'screen', 0.016133979), (
↪ 'n04141975', 'scale', 0.011708578), ('n03782006', 'monitor', 0.0091574294)]])
```

1.4 Advanced

The `Adversarial` class provides an advanced way to specify the adversarial example that should be found by an attack and provides detailed information about the created adversarial. In addition, it provides a way to improve a previously found adversarial example by re-running an attack.

1.4.1 Implicit

```
model = TensorFlowModel(images, logits, bounds=(0, 255))
criterion = TargetClassProbability('ostrich', p=0.99)
attack = LBFGSAttack(model, criterion)
```

Running the attack by passing image and label will implicitly create an `Adversarial` instance. By passing `unpack=False` we tell the attack to return the `Adversarial` instance rather than the actual image.

```
adversarial = attack(image, label=label, unpack=False)
```

We can then get the actual image using the `image` attribute:

```
adversarial_image = adversarial.image
```

1.4.2 Explicit

```
model = TensorFlowModel(images, logits, bounds=(0, 255))
criterion = TargetClassProbability('ostrich', p=0.99)
attack = LBFGSAttack()
```

We can also create the `Adversarial` instance ourselves and then pass it to the attack.

```
adversarial = Adversarial(model, criterion, image, label)
attack(adversarial)
```

Again, we can get the image using the `image` attribute:

```
adversarial_image = adversarial.image
```

This approach gives us more flexibility and allows us to specify a different distance measure:

```
distance = MeanAbsoluteDistance
adversarial = Adversarial(model, criterion, image, label, distance=distance)
```

1.5 Model Zoo

This tutorial will show you how the model zoo can be used to run your attack against a robust model.

1.5.1 Downloading a model

For this tutorial, we will download the *Analysis by Synthesis* model implemented in *PyTorch* and run a *FGSM (GradientSignAttack)* against it.

```
from foolbox import zoo

# download the model
model = zoo.get_model(url="https://github.com/bethgelab/AnalysisBySynthesis")

# read image and label
image = ...
label = ...

# apply attack on source image
attack = foolbox.attacks.FGSM(model)
adversarial = attack(image[:, :, :-1], label)
```

1.6 Development

To install Foolbox in editable mode, see the installation instructions under *Contributing to Foolbox*.

1.6.1 Running Tests

pytest

To run the tests, you need to have `pytest` and `pytest-cov` installed. Afterwards, you can simply run `pytest` in the root folder of the project. Some tests will require TensorFlow, PyTorch and the other frameworks, so to run all tests, you need to have all of them installed.

flake8

Foolbox follows the [PEP 8 style guide for Python code](#). To check for violations, we use `flake8` and run it like this:

```
flake8 --ignore E402,E741 .
```

1.6.2 New Adversarial Attacks

Foolbox makes it easy to develop new adversarial attacks that can be applied to arbitrary models.

To implement an attack, simply subclass the `Attack` class, implement the `__call__()` method and decorate it with the `:decorator:'call_decorator'`. The `:decorator:'call_decorator'` will make sure that your `__call__()` implementation will be called with an instance of the `Adversarial` class. You can use this instance to ask for model predictions and gradients, get the original image and its label and more. In addition, the `Adversarial` instance automatically keeps track of the best adversarial amongst all the images tested by the attack. That way, the implementation of the attack can focus on the attack logic.

1.7 FAQ

How does Foolbox handle inputs that are misclassified without any perturbation? The attacks will not be run and instead the unperturbed input is returned as an *adversarial* with distance 0 to the clean input.

What happens if an attack fails? The attack will return *None* and the distance will be *np.inf*.

Why is the returned adversarial not misclassified by my model? Most likely you have a discrepancy between how you evaluate your model and how you told Foolbox to evaluate it. For example, you might not be using the same preprocessing. Compare the output of the *predictions* method of the Foolbox model instance with your model's output (logits). This problem can also be caused by non-deterministic models. Make sure that your model is not stochastic and always returns the same output when given the same input. In rare cases it can also be that a seemingly deterministic model becomes numerically stochastic around the decision boundary (e.g. because of non-deterministic floating point *reduce_sum* operations). You can always check *adversarial.output* and *adversarial.adversarial_class* to see the output Foolbox got from your model when deciding that this was an adversarial.

Why are the gradients multiplied by the bounds (*max_ - min_*)? This scaling is meant to make hyperparameters such as the *epsilon* for FGSM independent of the bounds. *epsilon = 0.1* thus means that you perturb the image by 10% relative to the *max - min* range (which could for example go from 0 to 1 or from 0 to 255).

1.8 foolbox.models

Provides classes to wrap existing models in different frameworks so that they provide a unified API to the attacks.

1.8.1 Models

<code>Model</code>	Base class to provide attacks with a unified interface to models.
<code>DifferentiableModel</code>	Base class for differentiable models that provide gradients.

Continued on next page

Table 1 – continued from previous page

<i>TensorFlowModel</i>	Creates a <i>Model</i> instance from existing <i>TensorFlow</i> tensors.
<i>TensorFlowEagerModel</i>	Creates a <i>Model</i> instance from a <i>TensorFlow</i> model using eager execution.
<i>PyTorchModel</i>	Creates a <i>Model</i> instance from a <i>PyTorch</i> module.
<i>KerasModel</i>	Creates a <i>Model</i> instance from a <i>Keras</i> model.
<i>TheanoModel</i>	Creates a <i>Model</i> instance from existing <i>Theano</i> tensors.
<i>LasagneModel</i>	Creates a <i>Model</i> instance from a <i>Lasagne</i> network.
<i>MXNetModel</i>	Creates a <i>Model</i> instance from existing <i>MXNet</i> symbols and weights.
<i>MXNetGluonModel</i>	Creates a <i>Model</i> instance from an existing <i>MXNet Gluon</i> Block.

1.8.2 Wrappers

<i>ModelWrapper</i>	Base class for models that wrap other models.
<i>DifferentiableModelWrapper</i>	Base class for models that wrap other models and provide gradient methods.
<i>ModelWithoutGradients</i>	Turns a model into a model without gradients.
<i>ModelWithEstimatedGradients</i>	Turns a model into a model with gradients estimated by the given gradient estimator.
<i>CompositeModel</i>	Combines predictions of a (black-box) model with the gradient of a (substitute) model.

1.8.3 Detailed description

class foolbox.models.**Model** (*bounds*, *channel_axis*, *preprocessing*=(0, 1))

Base class to provide attacks with a unified interface to models.

The *Model* class represents a model and provides a unified interface to its predictions. Subclasses must implement `batch_predictions` and `num_classes`.

Model instances can be used as context managers and subclasses can require this to allocate and release resources.

Parameters

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

predictions()

num_classes()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions (*image*)

Convenience method that calculates predictions for a single image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

'numpy.ndarray' Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

See also:

batch_predictions()

class foolbox.models.**DifferentiableModel** (*bounds, channel_axis, preprocessing=(0, 1)*)

Base class for differentiable models that provide gradients.

The *DifferentiableModel* class can be used as a base class for models that provide gradients. Subclasses must implement *predictions_and_gradient*.

A model should be considered differentiable based on whether it provides a *predictions_and_gradient()* method and a *gradient()* method, not based on whether it subclasses *DifferentiableModel*.

A differentiable model does not necessarily provide reasonable values for the gradients, the gradient can be wrong. It only guarantees that the relevant methods can be called.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

gradient()

gradient (*image, label*)

Calculates the gradient of the cross-entropy loss w.r.t. the image.

The default implementation calls `predictions_and_gradient`. Subclasses can provide more efficient implementations that only calculate the gradient.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

class `foolbox.models.TensorFlowModel` (*images, logits, bounds, channel_axis=3, preprocessing=(0, 1)*)

Creates a *Model* instance from existing *TensorFlow* tensors.

Parameters

images [*tensorflow.Tensor*] The input to the model, usually a *tensorflow.placeholder*.

logits [*tensorflow.Tensor*] The predictions of the model, before the softmax.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t. to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

classmethod from_keras (*model*, *bounds*, *input_shape=None*, *channel_axis=3*, *preprocessing=(0, 1)*)

Alternative constructor for a TensorFlowModel that accepts a *tf.keras.Model* instance.

Parameters

model [*tensorflow.keras.Model*] A *tensorflow.keras.Model* that accepts a single input tensor and returns a single output tensor representing logits.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

input_shape [tuple] The shape of a single input, e.g. (28, 28, 1) for MNIST. If None, tries to get the the shape from the model's `input_shape` attribute.

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

gradient (*image*, *label*)

Calculates the gradient of the cross-entropy loss w.r.t. the image.

The default implementation calls `predictions_and_gradient`. Subclasses can provide more efficient implementations that only calculate the gradient.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

num_classes()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

class foolbox.models.**TensorFlowEagerModel** (*model, bounds, num_classes=None, channel_axis=3, preprocessing=(0, 1)*)

Creates a *Model* instance from a *TensorFlow* model using eager execution.

Parameters

model [a TensorFlow eager model] The TensorFlow eager model that should be attacked. It will be called with input tensors and should return logits.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

num_classes [int] If None, will try to infer it from the model's output shape.

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwises preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [*int*] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

`gradient()`

class `foolbox.models.PyTorchModel` (*model, bounds, num_classes, channel_axis=1, device=None, preprocessing=(0, 1)*)

Creates a *Model* instance from a *PyTorch* module.

Parameters

model [*torch.nn.Module*] The PyTorch model that should be attacked.

bounds [*tuple*] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

num_classes [*int*] Number of classes for which the model will output predictions.

channel_axis [*int*] The index of the axis that represents color channels.

device [*string*] A string specifying the device to do computation on. If None, will default to "cuda:0" if `torch.cuda.is_available()` or "cpu" if not.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient*, *image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image*, *label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

`gradient()`

class foolbox.models.**KerasModel** (*model*, *bounds*, *channel_axis=3*, *preprocessing=(0, 1)*, *predicts='probabilities'*)

Creates a *Model* instance from a *Keras* model.

Parameters

model [*keras.models.Model*] The *Keras* model that should be attacked.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

predicts [str] Specifies whether the *Keras* model predicts logits or probabilities. Logits are preferred, but probabilities are the default.

backward (*gradient*, *image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image*, *label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

`gradient()`

class `foolbox.models.TheanoModel` (*images, logits, bounds, num_classes, channel_axis=1, preprocessing=[0, 1]*)

Creates a *Model* instance from existing *Theano* tensors.

Parameters

images [*theano.tensor*] The input to the model.

logits [*theano.tensor*] The predictions of the model, before the softmax.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

num_classes [int] Number of classes for which the model will output predictions.

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwises preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

gradient (*image, label*)

Calculates the gradient of the cross-entropy loss w.r.t. the image.

The default implementation calls `predictions_and_gradient`. Subclasses can provide more efficient implementations that only calculate the gradient.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [*int*] Reference label used to calculate the gradient.

Returns

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient ()

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [*int*] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient ()

class `foolbox.models.LasagneModel` (*input_layer, logits_layer, bounds, channel_axis=1, preprocessing=(0, 1)*)

Creates a *Model* instance from a *Lasagne* network.

Parameters

input_layer [*lasagne.layers.Layer*] The input to the model.

logits_layer [*lasagne.layers.Layer*] The output of the model, before the softmax.

bounds [*tuple*] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [*int*] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwises preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient*, *image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

gradient ()

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

predictions ()

gradient (*image*, *label*)

Calculates the gradient of the cross-entropy loss w.r.t. the image.

The default implementation calls `predictions_and_gradient`. Subclasses can provide more efficient implementations that only calculate the gradient.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [*int*] Reference label used to calculate the gradient.

Returns

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient ()

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image*, *label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

class `foolbox.models.MXNetModel` (*data, logits, args, ctx, num_classes, bounds, channel_axis=1, aux_states=None, preprocessing=(0, 1)*)

Creates a *Model* instance from existing *MXNet* symbols and weights.

Parameters

data [*mxnet.symbol.Variable*] The input to the model.

logits [*mxnet.symbol.Symbol*] The predictions of the model, before the softmax.

args [*dictionary mapping str to mxnet.nd.array*] The parameters of the model.

ctx [*mxnet.context.Context*] The device, e.g. `mxnet.cpu()` or `mxnet.gpu()`.

num_classes [int] The number of classes.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [int] The index of the axis that represents color channels.

aux_states [*dictionary mapping str to mxnet.nd.array*] The states of auxiliary parameters of the model.

preprocessing: 2-element tuple with floats or numpy arrays Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t. to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

gradient()

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

`gradient()`

class foolbox.models.**MXNetGluonModel** (*block, bounds, num_classes, ctx=None, channel_axis=1, preprocessing=(0, 1)*)

Creates a *Model* instance from an existing *MXNet Gluon* Block.

Parameters

block [*mxnet.gluon.Block*] The Gluon Block representing the model to be run.

ctx [*mxnet.context.Context*] The device, e.g. `mxnet.cpu()` or `mxnet.gpu()`.

num_classes [int] The number of classes.

bounds [tuple] Tuple of lower and upper bound for the pixel values, usually (0, 1) or (0, 255).

channel_axis [int] The index of the axis that represents color channels.

preprocessing: 2-element tuple with floats or numpy arrays Elementwises preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

backward (*gradient, image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

`predictions()`

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image, label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

`gradient()`

class `foolbox.models.ModelWrapper` (*model*)

Base class for models that wrap other models.

This base class can be used to implement model wrappers that turn models into new models, for example by preprocessing the input or modifying the gradient.

Parameters

model [*Model*] The model that is wrapped.

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

predictions()

num_classes()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions (*image*)

Convenience method that calculates predictions for a single image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

'numpy.ndarray' Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

See also:

batch_predictions()

class foolbox.models.**DifferentiableModelWrapper** (*model*)

Base class for models that wrap other models and provide gradient methods.

This base class can be used to implement model wrappers that turn models into new models, for example by preprocessing the input or modifying the gradient.

Parameters

model [*Model*] The model that is wrapped.

class foolbox.models.**ModelWithoutGradients** (*model*)

Turns a model into a model without gradients.

class foolbox.models.**ModelWithEstimatedGradients** (*model*, *gradient_estimator*)

Turns a model into a model with gradients estimated by the given gradient estimator.

Parameters

model [*Model*] The model that is wrapped.

gradient_estimator [*callable*] Callable taking three arguments (pred_fn, image, label) and returning the estimated gradients. pred_fn will be the batch_predictions method of the wrapped model.

class foolbox.models.**CompositeModel** (*forward_model*, *backward_model*)

Combines predictions of a (black-box) model with the gradient of a (substitute) model.

Parameters

forward_model [*Model*] The model that should be fooled and will be used for predictions.

backward_model [*Model*] The model that provides the gradients.

backward (*gradient*, *image*)

Backpropagates the gradient of some loss w.r.t. the logits through the network and returns the gradient of that loss w.r.t to the input image.

Parameters

gradient [*numpy.ndarray*] Gradient of some loss w.r.t. the logits.

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [*numpy.ndarray*] The gradient w.r.t the image.

See also:

gradient ()

batch_predictions (*images*)

Calculates predictions for a batch of images.

Parameters

images [*numpy.ndarray*] Batch of inputs with shape as expected by the model.

Returns

'numpy.ndarray' Predictions (logits, i.e. before the softmax) with shape (batch size, number of classes).

See also:

predictions ()

gradient (*image*, *label*)

Calculates the gradient of the cross-entropy loss w.r.t. the image.

The default implementation calls `predictions_and_gradient`. Subclasses can provide more efficient implementations that only calculate the gradient.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [*int*] Reference label used to calculate the gradient.

Returns

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient ()

num_classes ()

Determines the number of classes.

Returns

int The number of classes for which the model creates predictions.

predictions_and_gradient (*image*, *label*)

Calculates predictions for an image and the gradient of the cross-entropy loss w.r.t. the image.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

label [int] Reference label used to calculate the gradient.

Returns

predictions [*numpy.ndarray*] Vector of predictions (logits, i.e. before the softmax) with shape (number of classes,).

gradient [*numpy.ndarray*] The gradient of the cross-entropy loss w.r.t. the image. Will have the same shape as the image.

See also:

gradient()

1.9 foolbox.criteria

Provides classes that define what is adversarial.

1.9.1 Criteria

We provide criteria for untargeted and targeted adversarial attacks.

<i>Misclassification</i>	Defines adversarials as images for which the predicted class is not the original class.
<i>TopKMisclassification</i>	Defines adversarials as images for which the original class is not one of the top k predicted classes.
<i>OriginalClassProbability</i>	Defines adversarials as images for which the probability of the original class is below a given threshold.
<i>ConfidentMisclassification</i>	Defines adversarials as images for which the probability of any class other than the original is above a given threshold.
<i>TargetClass</i>	Defines adversarials as images for which the predicted class is the given target class.
<i>TargetClassProbability</i>	Defines adversarials as images for which the probability of a given target class is above a given threshold.

1.9.2 Examples

Untargeted criteria:

```
>>> from foolbox.criteria import Misclassification
>>> criterion1 = Misclassification()
```

```
>>> from foolbox.criteria import TopKMisclassification
>>> criterion2 = TopKMisclassification(k=5)
```


Targeted criteria:

```
>>> from foolbox.criteria import TargetClass
>>> criterion3 = TargetClass(22)
```

```
>>> from foolbox.criteria import TargetClassProbability
>>> criterion4 = TargetClassProbability(22, p=0.99)
```

Criteria can be combined to create a new criterion:

```
>>> criterion5 = criterion2 & criterion3
```

1.9.3 Detailed description

class foolbox.criteria.**Criterion**

Base class for criteria that define what is adversarial.

The *Criterion* class represents a criterion used to determine if predictions for an image are adversarial given a reference label. It should be subclassed when implementing new criteria. Subclasses must implement `is_adversarial`.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [numpy.ndarray] A vector with the pre-softmax predictions for some image.

label [int] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by label, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class foolbox.criteria.**Misclassification**

Defines adversarials as images for which the predicted class is not the original class.

See also:

TopKMisclassification

Notes

Uses `numpy.argmax` to break ties.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [`numpy.ndarray`] A vector with the pre-softmax predictions for some image.

label [`int`] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by label, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class `foolbox.criteria.ConfidentMisclassification` (*p*)

Defines adversarials as images for which the probability of any class other than the original is above a given threshold.

Parameters

p [`float`] The threshold probability. If the probability of any class other than the original is at least *p*, the image is considered an adversarial. It must satisfy $0 \leq p \leq 1$.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [`numpy.ndarray`] A vector with the pre-softmax predictions for some image.

label [`int`] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by label, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class `foolbox.criteria.TopKMisclassification` (*k*)

Defines adversarials as images for which the original class is not one of the top *k* predicted classes.

For *k* = 1, the `Misclassification` class provides a more efficient implementation.

Parameters

k [int] Number of top predictions to which the reference label is compared to.

See also:

`Misclassification` Provides a more efficient implementation for *k* = 1.

Notes

Uses `numpy.argsort` to break ties.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [`numpy.ndarray`] A vector with the pre-softmax predictions for some image.

label [int] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by *label*, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class `foolbox.criteria.TargetClass` (*target_class*)

Defines adversarials as images for which the predicted class is the given target class.

Parameters

target_class [int] The target class that needs to be predicted for an image to be considered an adversarial.

Notes

Uses `numpy.argmax` to break ties.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [`numpy.ndarray`] A vector with the pre-softmax predictions for some image.

label [`int`] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by label, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class `foolbox.criteria.OriginalClassProbability` (*p*)

Defines adversarials as images for which the probability of the original class is below a given threshold.

This criterion alone does not guarantee that the class predicted for the adversarial image is not the original class (unless $p < 1 / \text{number of classes}$). Therefore, it should usually be combined with a classification criterion.

Parameters

p [`float`] The threshold probability. If the probability of the original class is below this threshold, the image is considered an adversarial. It must satisfy $0 \leq p \leq 1$.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [`numpy.ndarray`] A vector with the pre-softmax predictions for some image.

label [`int`] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by label, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

class `foolbox.criteria.TargetClassProbability` (*target_class, p*)

Defines adversarials as images for which the probability of a given target class is above a given threshold.

If the threshold is below 0.5, this criterion does not guarantee that the class predicted for the adversarial image is not the original class. In that case, it should usually be combined with a classification criterion.

Parameters

target_class [int] The target class for which the predicted probability must be above the threshold probability *p*, otherwise the image is not considered an adversarial.

p [float] The threshold probability. If the probability of the target class is above this threshold, the image is considered an adversarial. It must satisfy $0 \leq p \leq 1$.

is_adversarial (*predictions, label*)

Decides if predictions for an image are adversarial given a reference label.

Parameters

predictions [numpy.ndarray] A vector with the pre-softmax predictions for some image.

label [int] The label of the unperturbed reference image.

Returns

bool True if an image with the given predictions is an adversarial example when the ground-truth class is given by *label*, False otherwise.

name ()

Returns a human readable name that uniquely identifies the criterion with its hyperparameters.

Returns

str Human readable name that uniquely identifies the criterion with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

1.10 foolbox.zoo

1.10.1 Get Model

`foolbox.zoo.get_model` (*url, module_name='foolbox_model', **kwargs*)

Provides utilities to download foolbox-compatible robust models to easily test attacks against them by simply providing a git-URL.

Examples

Instantiate a model:

```
>>> from foolbox import zoo
>>> url = "https://github.com/bveliqi/foolbox-zoo-dummy.git"
>>> model = zoo.get_model(url) # doctest: +SKIP
```

Only works with a foolbox-zoo compatible repository. I.e. models need to have a `foolbox_model.py` file with a `create()`-function, which returns a foolbox-wrapped model.

Using the `kwargs` parameter it is possible to input an arbitrary number of parameters to this methods call. These parameters are forwarded to the instantiated model.

Example repositories:

- <https://github.com/bethgelab/AnalysisBySynthesis>
- https://github.com/bethgelab/mnist_challenge
- https://github.com/bethgelab/cifar10_challenge
- https://github.com/bethgelab/convex_adversarial
- <https://github.com/wielandbrendel/logit-pairing-foolbox.git>
- <https://github.com/bethgelab/defensive-distillation.git>

Parameters

- `url` – URL to the git repository
- `module_name` – the name of the module to import
- `kwargs` – Optional set of parameters that will be used by the to be instantiated model.

Returns a foolbox-wrapped model instance

1.10.2 Fetch Weights

`foolbox.zoo.fetch_weights(weights_uri, unzip=False)`

Provides utilities to download and extract packages containing model weights when creating foolbox-zoo compatible repositories, if the weights are not part of the repository itself.

Examples

Download and unzip weights:

```
>>> from foolbox import zoo
>>> url = 'https://github.com/MadryLab/mnist_challenge_models/raw/master/secret.
↳zip' # noqa F501
>>> weights_path = zoo.fetch_weights(url, unzip=True)
```

Parameters

- `weights_uri` – the URI to fetch the weights from
- `unzip` – should be `True` if the file to be downloaded is a zipped package

Returns local path where the weights have been downloaded and potentially unzipped to

1.11 foolbox.distances

Provides classes to measure the distance between images.

1.11.1 Distances

<i>MeanSquaredDistance</i>	Calculates the mean squared error between two images.
<i>MeanAbsoluteDistance</i>	Calculates the mean absolute error between two images.
<i>Linfinity</i>	Calculates the L-infinity norm of the difference between two images.
<i>L0</i>	Calculates the L0 norm of the difference between two images.

1.11.2 Aliases

<i>MSE</i>	alias of <i>foolbox.distances.MeanSquaredDistance</i>
<i>MAE</i>	alias of <i>foolbox.distances.MeanAbsoluteDistance</i>
<i>Linf</i>	alias of <i>foolbox.distances.Linfinity</i>

1.11.3 Base class

To implement a new distance, simply subclass the *Distance* class and implement the `_calculate()` method.

<i>Distance</i>	Base class for distances.
-----------------	---------------------------

1.11.4 Detailed description

class `foolbox.distances.Distance` (*reference=None, other=None, bounds=None, value=None*)
Base class for distances.

This class should be subclassed when implementing new distances. Subclasses must implement `_calculate`.

class `foolbox.distances.MeanSquaredDistance` (*reference=None, other=None, bounds=None, value=None*)
Calculates the mean squared error between two images.

class `foolbox.distances.MeanAbsoluteDistance` (*reference=None, other=None, bounds=None, value=None*)
Calculates the mean absolute error between two images.

class `foolbox.distances.Linfinity` (*reference=None, other=None, bounds=None, value=None*)
Calculates the L-infinity norm of the difference between two images.

class `foolbox.distances.L0` (*reference=None, other=None, bounds=None, value=None*)
Calculates the L0 norm of the difference between two images.

`foolbox.distances.MSE`
alias of `foolbox.distances.MeanSquaredDistance`

`foolbox.distances.MAE`
alias of `foolbox.distances.MeanAbsoluteDistance`

`foolbox.distances.Linf`
alias of `foolbox.distances.Linfinity`

1.12 foolbox.attacks

1.12.1 Gradient-based attacks

```
class foolbox.attacks.GradientAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                     object>, distance=<class 'fool-
                                     box.distances.MeanSquaredDistance'>, thresh-
                                     old=None)
```

Perturbs the image with the gradient of the loss w.r.t. the image, gradually increasing the magnitude until the image is misclassified.

Does not do anything if the model does not have a gradient.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000, max_epsilon=1)
```

Perturbs the image with the gradient of the loss w.r.t. the image, gradually increasing the magnitude until the image is misclassified.

Parameters

input_or_adv [`numpy.ndarray` or `Adversarial`] The original, unperturbed input as a `numpy.ndarray` or an `Adversarial` instance.

label [`int`] The reference label of the original input. Must be passed if `a` is a `numpy.ndarray`, must not be passed if `a` is an `Adversarial` instance.

unpack [`bool`] If true, returns the adversarial input, otherwise returns the `Adversarial` object.

epsilons [`int` or `Iterable[float]`] Either `Iterable` of step sizes in the gradient direction or number of step sizes between 0 and `max_epsilon` that should be tried.

max_epsilon [`float`] Largest step size if `epsilons` is not an `iterable`.

```
class foolbox.attacks.GradientSignAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                          object>, distance=<class 'fool-
                                          box.distances.MeanSquaredDistance'>, thresh-
                                          old=None)
```

Adds the sign of the gradient to the image, gradually increasing the magnitude until the image is misclassified. This attack is often referred to as Fast Gradient Sign Method and was introduced in [1].

Does not do anything if the model does not have a gradient.

References

[1]

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000, max_epsilon=1)
```

Adds the sign of the gradient to the image, gradually increasing the magnitude until the image is misclassified.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or Iterable[float]] Either Iterable of step sizes in the direction of the sign of the gradient or number of step sizes between 0 and `max_epsilon` that should be tried.

max_epsilon [float] Largest step size if `epsilons` is not an iterable.

`foolbox.attacks.FGSM`

alias of `foolbox.attacks.gradient.GradientSignAttack`

```
class foolbox.attacks.LinfinityBasicIterativeAttack (model=None, criterion=<foolbox.criteria.Misclassification
object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
threshold=None)
```

The Basic Iterative Method introduced in [1].

This attack is also known as Projected Gradient Descent (PGD) (without random start) or FGSM^k.

References

See also:

[*ProjectedGradientDescentAttack*](#)

[1]

```
__call__ (input_or_adv, label=None, unpack=True, binary_search=True, epsilon=0.3, stepsize=0.05,
iterations=10, random_start=False, return_early=True)
```

Simple iterative gradient-based attack known as Basic Iterative Method, Projected Gradient Descent or FGSM^k.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search [bool or int] Whether to perform a binary search over `epsilon` and `stepsize`, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if `binary_search` is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if `binary_search` is True, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

`foolbox.attacks.BasicIterativeMethod`

alias of `foolbox.attacks.iterative_projected_gradient.LinfinityBasicIterativeAttack`

`foolbox.attacks.BIM`

alias of `foolbox.attacks.iterative_projected_gradient.LinfinityBasicIterativeAttack`

```
class foolbox.attacks.L1BasicIterativeAttack (model=None, criterion=<foolbox.criteria.Misclassification
object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
threshold=None)
```

Modified version of the Basic Iterative Method that minimizes the L1 distance.

See also:

[`LinfinityBasicIterativeAttack`](#)

```
__call__ (input_or_adv, label=None, unpack=True, binary_search=True, epsilon=0.3, stepsize=0.05,
iterations=10, random_start=False, return_early=True)
```

Simple iterative gradient-based attack known as Basic Iterative Method, Projected Gradient Descent or FGSM^k.

Parameters

input_or_adv [`numpy.ndarray` or `Adversarial`] The original, unperturbed input as a `numpy.ndarray` or an `Adversarial` instance.

label [int] The reference label of the original input. Must be passed if *a* is a `numpy.ndarray`, must not be passed if *a* is an `Adversarial` instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the `Adversarial` object.

binary_search [bool or int] Whether to perform a binary search over epsilon and stepsize, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if `binary_search` is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if `binary_search` is True, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

```
class foolbox.attacks.L2BasicIterativeAttack (model=None, criterion=<foolbox.criteria.Misclassification
object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
threshold=None)
```

Modified version of the Basic Iterative Method that minimizes the L2 distance.

See also:

LinfinitBasicIterativeAttack

__call__ (*input_or_adv*, *label=None*, *unpack=True*, *binary_search=True*, *epsilon=0.3*, *stepsize=0.05*, *iterations=10*, *random_start=False*, *return_early=True*)

Simple iterative gradient-based attack known as Basic Iterative Method, Projected Gradient Descent or FGSM^k.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search [bool or int] Whether to perform a binary search over epsilon and stepsize, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if *binary_search* is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if *binary_search* is True, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

```
class foolbox.attacks.ProjectedGradientDescentAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

The Projected Gradient Descent Attack introduced in [1] without random start.

When used without a random start, this attack is also known as Basic Iterative Method (BIM) or FGSM^k.

References

See also:

LinfinitBasicIterativeAttack and *RandomStartProjectedGradientDescentAttack*

[1]

__call__ (*input_or_adv*, *label=None*, *unpack=True*, *binary_search=True*, *epsilon=0.3*, *stepsize=0.01*, *iterations=40*, *random_start=False*, *return_early=True*)

Simple iterative gradient-based attack known as Basic Iterative Method, Projected Gradient Descent or FGSM^k.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search [bool or int] Whether to perform a binary search over epsilon and stepsize, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if *binary_search* is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if *binary_search* is True, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

`foolbox.attacks.ProjectedGradientDescent`

alias of `foolbox.attacks.iterative_projected_gradient.ProjectedGradientDescentAttack`

```
class foolbox.attacks.RandomStartProjectedGradientDescentAttack (model=None,
                                                                crite-
                                                                rion=<foolbox.criteria.Misclassification
                                                                object>, dis-
                                                                tance=<class
                                                                'fool-
                                                                box.distances.MeanSquaredDistance'>
                                                                thresh-
                                                                old=None)
```

The Projected Gradient Descent Attack introduced in [1] with random start.

References

See also:

ProjectedGradientDescentAttack

[1]

```
__call__ (input_or_adv, label=None, unpack=True, binary_search=True, epsilon=0.3, stepsize=0.01,
          iterations=40, random_start=True, return_early=True)
```

Simple iterative gradient-based attack known as Basic Iterative Method, Projected Gradient Descent or FGSM^k.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search [bool or int] Whether to perform a binary search over epsilon and stepsize, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if `binary_search` is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if `binary_search` is True, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

```
foolbox.attacks.RandomProjectedGradientDescent
    alias          of          foolbox.attacks.iterative_projected_gradient.
    RandomStartProjectedGradientDescentAttack
```

```
foolbox.attacks.RandomPGD
    alias          of          foolbox.attacks.iterative_projected_gradient.
    RandomStartProjectedGradientDescentAttack
```

```
class foolbox.attacks.MomentumIterativeAttack (model=None,          crite-
                                               ri-
                                               on=<foolbox.criteria.Misclassification
                                               object>,   distance=<class   'fool-
                                               box.distances.MeanSquaredDistance'>,
                                               threshold=None)
```

The Momentum Iterative Method attack introduced in [1]. It's like the Basic Iterative Method or Projected Gradient Descent except that it uses momentum.

References

[1]

```
__call__ (input_or_adv, label=None, unpack=True, binary_search=True, epsilon=0.3, stepsize=0.06,
         iterations=10, decay_factor=1.0, random_start=False, return_early=True)
    Momentum-based iterative gradient attack known as Momentum Iterative Method.
```

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search [bool] Whether to perform a binary search over epsilon and stepsize, keeping their ratio constant and using their values to start the search. If False, hyperparameters are not optimized. Can also be an integer, specifying the number of binary search steps (default 20).

epsilon [float] Limit on the perturbation size; if `binary_search` is True, this value is only for initialization and automatically adapted.

stepsize [float] Step size for gradient descent; if `binary_search` is `True`, this value is only for initialization and automatically adapted.

iterations [int] Number of iterations for each gradient descent run.

decay_factor [float] Decay factor used by the momentum term.

random_start [bool] Start the attack from a random point rather than from the original input.

return_early [bool] Whether an individual gradient descent run should stop as soon as an adversarial is found.

`foolbox.attacks.MomentumIterativeMethod`

alias of `foolbox.attacks.iterative_projected_gradient.MomentumIterativeAttack`

class `foolbox.attacks.LBFGSAttack` (**args*, ***kwargs*)

Uses L-BFGS-B to minimize the distance between the image and the adversarial as well as the cross-entropy between the predictions for the adversarial and the the one-hot encoded target class.

If the criterion does not have a target class, a random class is chosen from the set of all classes except the original one.

Notes

This implementation generalizes algorithm 1 in [1] to support other targeted criteria and other distance measures.

References

[1]

`__call__` (*input_or_adv*, *label=None*, *unpack=True*, *epsilon=1e-05*, *num_random_targets=0*, *max_iter=150*)

Uses L-BFGS-B to minimize the distance between the image and the adversarial as well as the cross-entropy between the predictions for the adversarial and the the one-hot encoded target class.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilon [float] Epsilon of the binary search.

num_random_targets [int] Number of random target classes if no target class is given by the criterion.

maxiter [int] Maximum number of iterations for L-BFGS-B.

`__init__` (**args*, ***kwargs*)

Initialize self. See `help(type(self))` for accurate signature.

`name` ()

Returns a human readable name that uniquely identifies the attack with its hyperparameters.

Returns

str Human readable name that uniquely identifies the attack with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

```
class foolbox.attacks.DeepFoolAttack (model=None, criterion=<foolbox.criteria.Misclassification
object>, distance=<class 'fool-
box.distances.MeanSquaredDistance'>, thresh-
old=None)
```

Simple and close to optimal gradient-based adversarial attack.

Implements DeepFool introduced in [1].

References

[1]

__call__ (*input_or_adv, label=None, unpack=True, steps=100, subsample=10, p=None*)

Simple and close to optimal gradient-based adversarial attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

steps [int] Maximum number of steps to perform.

subsample [int] Limit on the number of the most likely classes that should be considered. A small value is usually sufficient and much faster.

p [int or float] Lp-norm that should be minimized, must be 2 or *np.inf*.

```
class foolbox.attacks.NewtonFoolAttack (model=None, criterion=<foolbox.criteria.Misclassification
object>, distance=<class 'fool-
box.distances.MeanSquaredDistance'>, thresh-
old=None)
```

Implements the NewtonFool Attack.

The attack was introduced in [1].

References

[1]

__call__ (*input_or_adv, label=None, unpack=True, max_iter=100, eta=0.01*)

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

max_iter [int] The maximum number of iterations.

eta [float] the eta coefficient

```
class foolbox.attacks.DeepFoolL2Attack (model=None, criterion=<foolbox.criteria.Misclassification
                                     object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                     threshold=None)
```

__call__ (input_or_adv, label=None, unpack=True, steps=100, subsample=10)
Simple and close to optimal gradient-based adversarial attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

steps [int] Maximum number of steps to perform.

subsample [int] Limit on the number of the most likely classes that should be considered. A small value is usually sufficient and much faster.

p [int or float] Lp-norm that should be minimized, must be 2 or np.inf.

```
class foolbox.attacks.DeepFoolLinfinityAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                              object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                              threshold=None)
```

__call__ (input_or_adv, label=None, unpack=True, steps=100, subsample=10)
Simple and close to optimal gradient-based adversarial attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

steps [int] Maximum number of steps to perform.

subsample [int] Limit on the number of the most likely classes that should be considered. A small value is usually sufficient and much faster.

p [int or float] Lp-norm that should be minimized, must be 2 or np.inf.

```
class foolbox.attacks.ADefAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                threshold=None)
```

Adversarial attack that distorts the image, i.e. changes the locations of pixels. The algorithm is described in [Rf241e6d2664d-1], a Repository with the original code can be found in [Rf241e6d2664d-2]. References ——— .. [Rf241e6d2664d-1] Rima Alaifari, Giovanni S. Alberti, and Tandri Gauksson:

“ADef: an Iterative Algorithm to Construct Adversarial Deformations”, <https://arxiv.org/abs/1804.07729>

```
__call__(input_or_adv, unpack=True, max_iter=100, max_norm=<Mock name='mock.inf'
        id='140200757448320'>, label=None, smooth=1.0, subsample=10)
```

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

max_iter [int > 0] Maximum number of iterations (default max_iter = 100).

max_norm [float] Maximum l2 norm of vector field (default max_norm = numpy.inf).

smooth [float >= 0] Width of the Gaussian kernel used for smoothing. (default is smooth = 0 for no smoothing).

subsample [int >= 2] Limit on the number of the most likely classes that should be considered. A small value is usually sufficient and much faster. (default subsample = 10)

```
class foolbox.attacks.SLSQPAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                object>, distance=<class 'fool-
                                box.distances.MeanSquaredDistance'>, threshold=None)
```

Uses SLSQP to minimize the distance between the image and the adversarial under the constraint that the image is adversarial.

```
__call__(input_or_adv, label=None, unpack=True)
```

Uses SLSQP to minimize the distance between the image and the adversarial under the constraint that the image is adversarial.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, correctly classified image. If image is a *numpy* array, label must be passed as well. If image is an *Adversarial* instance, label must not be passed.

label [int] The reference label of the original image. Must be passed if image is a *numpy* array, must not be passed if image is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial image, otherwise returns the *Adversarial* object.

```
class foolbox.attacks.SaliencyMapAttack (model=None, criterion=<foolbox.criteria.Misclassification
                                        object>, distance=<class 'fool-
                                        box.distances.MeanSquaredDistance'>, threshold=None)
```

Implements the Saliency Map Attack.

The attack was introduced in [1].

References

[1]

```
__call__(input_or_adv, label=None, unpack=True, max_iter=2000, num_random_targets=0,
         fast=True, theta=0.1, max_perturbations_per_pixel=7)
```

Implements the Saliency Map Attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

max_iter [int] The maximum number of iterations to run.

num_random_targets [int] Number of random target classes if no target class is given by the criterion.

fast [bool] Whether to use the fast saliency map calculation.

theta [float] perturbation per pixel relative to [min, max] range.

max_perturbations_per_pixel [int] Maximum number of times a pixel can be modified.

```
class foolbox.attacks.IterativeGradientAttack(model=None, criterion=<foolbox.criteria.Misclassification
                                             object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                             threshold=None)
```

Like *GradientAttack* but with several steps for each epsilon.

```
__call__(input_or_adv, label=None, unpack=True, epsilons=100, max_epsilon=1, steps=10)
```

Like *GradientAttack* but with several steps for each epsilon.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or *Iterable*[float]] Either *Iterable* of step sizes in the gradient direction or number of step sizes between 0 and *max_epsilon* that should be tried.

max_epsilon [float] Largest step size if *epsilons* is not an *iterable*.

steps [int] Number of iterations to run.

```
class foolbox.attacks.IterativeGradientSignAttack(model=None, criterion=<foolbox.criteria.Misclassification
                                                object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                                threshold=None)
```

Like *GradientSignAttack* but with several steps for each epsilon.

```
__call__(input_or_adv, label=None, unpack=True, epsilons=100, max_epsilon=1, steps=10)
```

Like *GradientSignAttack* but with several steps for each epsilon.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or Iterable[float]] Either Iterable of step sizes in the direction of the sign of the gradient or number of step sizes between 0 and *max_epsilon* that should be tried.

max_epsilon [float] Largest step size if *epsilons* is not an iterable.

steps [int] Number of iterations to run.

```
class foolbox.attacks.CarliniWagnerL2Attack (model=None, crite-
rion=<foolbox.criteria.Misclassification
object>, distance=<class 'fool-
box.distances.MeanSquaredDistance'>,
threshold=None)
```

The L2 version of the Carlini & Wagner attack.

This attack is described in [1]. This implementation is based on the reference implementation by Carlini [2]. For bounds (0, 1), it differs from [2] because we normalize the squared L2 loss with the bounds.

References

[1], [2]

```
__call__ (input_or_adv, label=None, unpack=True, binary_search_steps=5, max_iterations=1000,
confidence=0, learning_rate=0.005, initial_const=0.01, abort_early=True)
```

The L2 version of the Carlini & Wagner attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

binary_search_steps [int] The number of steps for the binary search used to find the optimal tradeoff-constant between distance and confidence.

max_iterations [int] The maximum number of iterations. Larger values are more accurate; setting it too small will require a large learning rate and will produce poor results.

confidence [int or float] Confidence of adversarial examples: a higher value produces adversarials that are further away, but more strongly classified as adversarial.

learning_rate [float] The learning rate for the attack algorithm. Smaller values produce better results but take longer to converge.

initial_const [float] The initial tradeoff-constant to use to tune the relative importance of distance and confidence. If *binary_search_steps* is large, the initial constant is not important.

abort_early [bool] If True, Adam will be aborted if the loss hasn't decreased for some time (a tenth of *max_iterations*).

```
static best_other_class (logits, exclude)
```

Returns the index of the largest logit, ignoring the class that is passed as *exclude*.

classmethod `loss_function` (*const, a, x, logits, reconstructed_original, confidence, min_, max_*)
Returns the loss and the gradient of the loss w.r.t. *x*, assuming that `logits = model(x)`.

1.12.2 Score-based attacks

class `foolbox.attacks.SinglePixelAttack` (*model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None*)

Perturbs just a single pixel and sets it to the min or max.

`__call__` (*input_or_adv, label=None, unpack=True, max_pixels=1000*)

Perturbs just a single pixel and sets it to the min or max.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, correctly classified image. If image is a *numpy* array, *label* must be passed as well. If image is an *Adversarial* instance, *label* must not be passed.

label [*int*] The reference label of the original image. Must be passed if image is a *numpy* array, must not be passed if image is an *Adversarial* instance.

unpack [*bool*] If true, returns the adversarial image, otherwise returns the *Adversarial* object.

max_pixels [*int*] Maximum number of pixels to try.

class `foolbox.attacks.LocalSearchAttack` (*model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None*)

A black-box attack based on the idea of greedy local search.

This implementation is based on the algorithm in [1].

References

[1]

`__call__` (*input_or_adv, label=None, unpack=True, r=1.5, p=10.0, d=5, t=5, R=150*)

A black-box attack based on the idea of greedy local search.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, correctly classified image. If image is a *numpy* array, *label* must be passed as well. If image is an *Adversarial* instance, *label* must not be passed.

label [*int*] The reference label of the original image. Must be passed if image is a *numpy* array, must not be passed if image is an *Adversarial* instance.

unpack [*bool*] If true, returns the adversarial image, otherwise returns the *Adversarial* object.

r [*float*] Perturbation parameter that controls the cyclic perturbation; must be in [0, 2]

p [*float*] Perturbation parameter that controls the pixel sensitivity estimation

d [int] The half side length of the neighborhood square

t [int] The number of pixels perturbed at each round

R [int] An upper bound on the number of iterations

class foolbox.attacks.**ApproximateLBFGSAttack** (*args, **kwargs)

Same as *LBFGSAttack* with `approximate_gradient` set to `True`.

`__init__` (*args, **kwargs)

Initialize self. See `help(type(self))` for accurate signature.

1.12.3 Decision-based attacks

class foolbox.attacks.**BoundaryAttack** (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)

A powerful adversarial attack that requires neither gradients nor probabilities.

This is the reference implementation for the attack introduced in [1].

Notes

This implementation provides several advanced features:

- ability to continue previous attacks by passing an instance of the Adversarial class
- ability to pass an explicit starting point; especially to initialize a targeted attack
- ability to pass an alternative attack used for initialization
- fine-grained control over logging
- ability to specify the batch size
- optional automatic batch size tuning
- optional multithreading for random number generation
- optional multithreading for candidate point generation

References

[1]

`__call__` (input_or_adv, label=None, unpack=True, iterations=5000, max_directions=25, starting_point=None, initialization_attack=None, log_every_n_steps=1, spherical_step=0.01, source_step=0.01, step_adaptation=1.5, batch_size=1, tune_batch_size=True, threaded_rnd=True, threaded_gen=True, alternative_generator=False, internal_dtype=<Mock name='mock.float64' id='140200757680280'>, verbose=False)

Applies the Boundary Attack.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, correctly classified image. If image is a *numpy* array, label must be passed as well. If image is an *Adversarial* instance, label must not be passed.

label [int] The reference label of the original image. Must be passed if image is a *numpy* array, must not be passed if image is an *Adversarial* instance.

- unpack** [bool] If true, returns the adversarial image, otherwise returns the Adversarial object.
- iterations** [int] Maximum number of iterations to run. Might converge and stop before that.
- max_directions** [int] Maximum number of trials per iteration.
- starting_point** [*numpy.ndarray*] Adversarial input to use as a starting point, in particular for targeted attacks.
- initialization_attack** [*Attack*] Attack to use to find a starting point. Defaults to *BlendedUniformNoiseAttack*.
- log_every_n_steps** [int] Determines verbosity of the logging.
- spherical_step** [float] Initial step size for the orthogonal (spherical) step.
- source_step** [float] Initial step size for the step towards the target.
- step_adaptation** [float] Factor by which the step sizes are multiplied or divided.
- batch_size** [int] Batch size or initial batch size if *tune_batch_size* is True
- tune_batch_size** [bool] Whether or not the batch size should be automatically chosen between 1 and *max_directions*.
- threaded_rnd** [bool] Whether the random number generation should be multithreaded.
- threaded_gen** [bool] Whether the candidate point generation should be multithreaded.
- alternative_generator: bool** Whether an alternative implementation of the candidate generator should be used.
- internal_dtype** [*np.float32* or *np.float64*] Higher precision might be slower but is numerically more stable.
- verbose** [bool] Controls verbosity of the attack.

class `foolbox.attacks.SpatialAttack` (*model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None*)

Adversarially chosen rotations and translations [1].

This implementation is based on the reference implementation by Madry et al.: https://github.com/MadryLab/adversarial_spatial

References

[1]

`__call__` (*input_or_adv, label=None, unpack=True, do_rotations=True, do_translations=True, x_shift_limits=(-5, 5), y_shift_limits=(-5, 5), angular_limits=(-5, 5), granularity=10, random_sampling=False, abort_early=True*)

Adversarially chosen rotations and translations.

Parameters

- input_or_adv** [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.
- label** [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.
- unpack** [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

do_rotations [bool] If False no rotations will be applied to the image.

do_translations [bool] If False no translations will be applied to the image.

x_shift_limits [int or (int, int)] Limits for horizontal translations in pixels. If one integer is provided the limits will be (-x_shift_limits, x_shift_limits).

y_shift_limits [int or (int, int)] Limits for vertical translations in pixels. If one integer is provided the limits will be (-y_shift_limits, y_shift_limits).

angular_limits [int or (int, int)] Limits for rotations in degrees. If one integer is provided the limits will be [-angular_limits, angular_limits].

granularity [int] Density of sampling within limits for each dimension.

random_sampling [bool] If True we sample translations/rotations randomly within limits, otherwise we use a regular grid.

abort_early [bool] If True, the attack stops as soon as it finds an adversarial.

```
class foolbox.attacks.PointwiseAttack (model=None,                                crite-
                                       rion=<foolbox.criteria.Misclassification
                                       object>,                                distance=<class
                                       box.distances.MeanSquaredDistance'>,    thresh-
                                       old=None)
```

Starts with an adversarial and performs a binary search between the adversarial and the original for each dimension of the input individually.

References

[1]

```
__call__ (input_or_adv, label=None, unpack=True, starting_point=None, initializa-
         tion_attack=None)
```

Starts with an adversarial and performs a binary search between the adversarial and the original for each dimension of the input individually.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

starting_point [*numpy.ndarray*] *Adversarial* input to use as a starting point, in particular for targeted attacks.

initialization_attack [*Attack*] *Attack* to use to find a starting point. Defaults to *SaltAndPepperNoiseAttack*.

```
class foolbox.attacks.GaussianBlurAttack (model=None,                                crite-
                                       rion=<foolbox.criteria.Misclassification
                                       object>,                                distance=<class
                                       box.distances.MeanSquaredDistance'>,    thresh-
                                       old=None)
```

Blurs the image until it is misclassified.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000)
```

Blurs the image until it is misclassified.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or Iterable[float]] Either Iterable of standard deviations of the Gaussian blur or number of standard deviations between 0 and 1 that should be tried.

```
class foolbox.attacks.ContrastReductionAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

Reduces the contrast of the image until it is misclassified.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000)  
Reduces the contrast of the image until it is misclassified.
```

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or Iterable[float]] Either Iterable of contrast levels or number of contrast levels between 1 and 0 that should be tried. Epsilons are one minus the contrast level.

```
class foolbox.attacks.AdditiveUniformNoiseAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

Adds uniform noise to the image, gradually increasing the standard deviation until the image is misclassified.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000)  
Adds uniform or Gaussian noise to the image, gradually increasing the standard deviation until the image is misclassified.
```

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or Iterable[float]] Either Iterable of noise levels or number of noise levels between 0 and 1 that should be tried.

```
__class__  
alias of abc.ABCMeta
```


`__delattr__` (*\$self, name, /*)
 Implement `delattr(self, name)`.

`__dir__` () → list
 default `dir()` implementation

`__eq__` (*\$self, value, /*)
 Return `self==value`.

`__format__` ()
 default object formatter

`__ge__` (*\$self, value, /*)
 Return `self>=value`.

`__getattr__` (*\$self, name, /*)
 Return `getattr(self, name)`.

`__gt__` (*\$self, value, /*)
 Return `self>value`.

`__hash__` (*\$self, /*)
 Return `hash(self)`.

`__init__` (*model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None*)
 Initialize self. See `help(type(self))` for accurate signature.

`__le__` (*\$self, value, /*)
 Return `self<=value`.

`__lt__` (*\$self, value, /*)
 Return `self<value`.

`__ne__` (*\$self, value, /*)
 Return `self!=value`.

`__new__` (*\$type, *args, **kwargs*)
 Create and return a new object. See `help(type)` for accurate signature.

`__reduce__` ()
 helper for pickle

`__reduce_ex__` ()
 helper for pickle

`__repr__` (*\$self, /*)
 Return `repr(self)`.

`__setattr__` (*\$self, name, value, /*)
 Implement `setattr(self, name, value)`.

`__sizeof__` () → int
 size of object in memory, in bytes

`__str__` (*\$self, /*)
 Return `str(self)`.

`__subclasshook__` ()
 Abstract classes can override this to customize `issubclass()`.

This is invoked early on by `abc.ABCMeta.__subclasscheck__()`. It should return `True`, `False` or `NotImplemented`. If it returns `NotImplemented`, the normal algorithm is used. Otherwise, it overrides the normal algorithm (and the outcome is cached).

`__weakref__`

list of weak references to the object (if defined)

`name()`

Returns a human readable name that uniquely identifies the attack with its hyperparameters.

Returns

str Human readable name that uniquely identifies the attack with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

```
class foolbox.attacks.AdditiveGaussianNoiseAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

Adds Gaussian noise to the image, gradually increasing the standard deviation until the image is misclassified.

`__call__` (*input_or_adv*, *label=None*, *unpack=True*, *epsilons=1000*)

Adds uniform or Gaussian noise to the image, gradually increasing the standard deviation until the image is misclassified.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or *Iterable*[float]] Either *Iterable* of noise levels or number of noise levels between 0 and 1 that should be tried.

`__class__`

alias of `abc.ABCMeta`

`__delattr__` (*\$self*, *name*, /)

Implement `delattr(self, name)`.

`__dir__` () → list

default `dir()` implementation

`__eq__` (*\$self*, *value*, /)

Return `self==value`.

`__format__` ()

default object formatter

`__ge__` (*\$self*, *value*, /)

Return `self>=value`.

`__getattr__` (*\$self*, *name*, /)

Return `getattr(self, name)`.

`__gt__` (*\$self*, *value*, /)

Return `self>value`.

__hash__ (*\$self, /*)
Return hash(self).

__init__ (*model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None*)
Initialize self. See help(type(self)) for accurate signature.

__le__ (*\$self, value, /*)
Return self<=value.

__lt__ (*\$self, value, /*)
Return self<value.

__ne__ (*\$self, value, /*)
Return self!=value.

__new__ (*\$type, *args, **kwargs*)
Create and return a new object. See help(type) for accurate signature.

__reduce__ ()
helper for pickle

__reduce_ex__ ()
helper for pickle

__repr__ (*\$self, /*)
Return repr(self).

__setattr__ (*\$self, name, value, /*)
Implement setattr(self, name, value).

__sizeof__ () → int
size of object in memory, in bytes

__str__ (*\$self, /*)
Return str(self).

__subclasshook__ ()
Abstract classes can override this to customize issubclass().

This is invoked early on by abc.ABCMeta.__subclasscheck__(). It should return True, False or NotImplemented. If it returns NotImplemented, the normal algorithm is used. Otherwise, it overrides the normal algorithm (and the outcome is cached).

__weakref__
list of weak references to the object (if defined)

name ()
Returns a human readable name that uniquely identifies the attack with its hyperparameters.

Returns

str Human readable name that uniquely identifies the attack with its hyperparameters.

Notes

Defaults to the class name but subclasses can provide more descriptive names and must take hyperparameters into account.

```
class foolbox.attacks.SaltAndPepperNoiseAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

Increases the amount of salt and pepper noise until the image is misclassified.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=100, repetitions=10)
```

Increases the amount of salt and pepper noise until the image is misclassified.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int] Number of steps to try between probability 0 and 1.

repetitions [int] Specifies how often the attack will be repeated.

```
class foolbox.attacks.BlendedUniformNoiseAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

Blends the image with a uniform noise image until it is misclassified.

```
__call__ (input_or_adv, label=None, unpack=True, epsilons=1000, max_directions=1000)
```

Blends the image with a uniform noise image until it is misclassified.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [int] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [bool] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

epsilons [int or *Iterable*[float]] Either *Iterable* of blending steps or number of blending steps between 0 and 1 that should be tried.

max_directions [int] Maximum number of random images to try.

1.12.4 Other attacks

```
class foolbox.attacks.BinarizationRefinementAttack (model=None, criterion=<foolbox.criteria.Misclassification object>, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None)
```

For models that preprocess their inputs by binarizing the inputs, this attack can improve adversarials found by other attacks. It does so by utilizing information about the binarization and mapping values to the corresponding value in the clean input or to the right side of the threshold.

`__call__` (*input_or_adv*, *label=None*, *unpack=True*, *starting_point=None*, *threshold=None*, *included_in='upper'*)

For models that preprocess their inputs by binarizing the inputs, this attack can improve adversarials found by other attacks. It does so by utilizing information about the binarization and mapping values to the corresponding value in the clean input or to the right side of the threshold.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [*int*] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [*bool*] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

starting_point [*numpy.ndarray*] Adversarial input to use as a starting point.

threshold [*float*] The threshold used by the models binarization. If none, defaults to $(\text{model.bounds}()[1] - \text{model.bounds}()[0]) / 2$.

included_in [*str*] Whether the threshold value itself belongs to the lower or upper interval.

class `foolbox.attacks.PrecomputedImagesAttack` (*input_images*, *output_images*, **args*, ***kwargs*)

Attacks a model using precomputed adversarial candidates.

Parameters

input_images [*numpy.ndarray*] The original images that will be expected by this attack.

output_images [*numpy.ndarray*] The adversarial candidates corresponding to the *input_images*.

***args** [*positional args*] Positional args passed to the *Attack* base class.

****kwargs** [*keyword args*] Keyword args passed to the *Attack* base class.

`__call__` (*input_or_adv*, *label=None*, *unpack=True*)

Attacks a model using precomputed adversarial candidates.

Parameters

input_or_adv [*numpy.ndarray* or *Adversarial*] The original, unperturbed input as a *numpy.ndarray* or an *Adversarial* instance.

label [*int*] The reference label of the original input. Must be passed if *a* is a *numpy.ndarray*, must not be passed if *a* is an *Adversarial* instance.

unpack [*bool*] If true, returns the adversarial input, otherwise returns the *Adversarial* object.

`__init__` (*input_images*, *output_images*, **args*, ***kwargs*)

Initialize self. See `help(type(self))` for accurate signature.

Gradient-based attacks

<i>GradientAttack</i>	Perturbs the image with the gradient of the loss w.r.t.
<i>GradientSignAttack</i>	Adds the sign of the gradient to the image, gradually increasing the magnitude until the image is misclassified.
<i>FGSM</i>	alias of <code>foolbox.attacks.gradient.GradientSignAttack</code>

Continued on next page

Table 8 – continued from previous page

<i>LinfinitBasicIterativeAttack</i>	The Basic Iterative Method introduced in [1].
<i>BasicIterativeMethod</i>	alias of foolbox.attacks.iterative_projected_gradient. LinfinitBasicIterativeAttack
<i>BIM</i>	alias of foolbox.attacks.iterative_projected_gradient. LinfinitBasicIterativeAttack
<i>L1BasicIterativeAttack</i>	Modified version of the Basic Iterative Method that minimizes the L1 distance.
<i>L2BasicIterativeAttack</i>	Modified version of the Basic Iterative Method that minimizes the L2 distance.
<i>ProjectedGradientDescentAttack</i>	The Projected Gradient Descent Attack introduced in [1] without random start.
<i>ProjectedGradientDescent</i>	alias of foolbox.attacks.iterative_projected_gradient. ProjectedGradientDescentAttack
<i>PGD</i>	alias of foolbox.attacks.iterative_projected_gradient. ProjectedGradientDescentAttack
<i>RandomStartProjectedGradientDescentAttack</i>	The Projected Gradient Descent Attack introduced in [1] with random start.
<i>RandomProjectedGradientDescent</i>	alias of foolbox.attacks.iterative_projected_gradient. RandomStartProjectedGradientDescentAttack
<i>RandomPGD</i>	alias of foolbox.attacks.iterative_projected_gradient. RandomStartProjectedGradientDescentAttack
<i>MomentumIterativeAttack</i>	The Momentum Iterative Method attack introduced in [1].
<i>MomentumIterativeMethod</i>	alias of foolbox.attacks.iterative_projected_gradient. MomentumIterativeAttack
<i>LBFSGSAttack</i>	Uses L-BFGS-B to minimize the distance between the image and the adversarial as well as the cross-entropy between the predictions for the adversarial and the the one-hot encoded target class.
<i>DeepFoolAttack</i>	Simple and close to optimal gradient-based adversarial attack.
<i>NewtonFoolAttack</i>	Implements the NewtonFool Attack.
<i>DeepFoolL2Attack</i>	
<i>DeepFoolLinfinitAttack</i>	
<i>ADefAttack</i>	Adversarial attack that distorts the image, i.e.
<i>SLSQPAttack</i>	Uses SLSQP to minimize the distance between the image and the adversarial under the constraint that the image is adversarial.
<i>SaliencyMapAttack</i>	Implements the Saliency Map Attack.
<i>IterativeGradientAttack</i>	Like GradientAttack but with several steps for each epsilon.
<i>IterativeGradientSignAttack</i>	Like GradientSignAttack but with several steps for each epsilon.
<i>CarliniWagnerL2Attack</i>	The L2 version of the Carlini & Wagner attack.

Score-based attacks

<i>SinglePixelAttack</i>	Perturbs just a single pixel and sets it to the min or max.
<i>LocalSearchAttack</i>	A black-box attack based on the idea of greedy local search.
<i>ApproximateLBFGSAttack</i>	Same as <i>LBFGSAttack</i> with <code>approximate_gradient</code> set to <code>True</code> .

Decision-based attacks

<i>BoundaryAttack</i>	A powerful adversarial attack that requires neither gradients nor probabilities.
<i>SpatialAttack</i>	Adversarially chosen rotations and translations [1].
<i>PointwiseAttack</i>	Starts with an adversarial and performs a binary search between the adversarial and the original for each dimension of the input individually.
<i>GaussianBlurAttack</i>	Blurs the image until it is misclassified.
<i>ContrastReductionAttack</i>	Reduces the contrast of the image until it is misclassified.
<i>AdditiveUniformNoiseAttack</i>	Adds uniform noise to the image, gradually increasing the standard deviation until the image is misclassified.
<i>AdditiveGaussianNoiseAttack</i>	Adds Gaussian noise to the image, gradually increasing the standard deviation until the image is misclassified.
<i>SaltAndPepperNoiseAttack</i>	Increases the amount of salt and pepper noise until the image is misclassified.
<i>BlendedUniformNoiseAttack</i>	Blends the image with a uniform noise image until it is misclassified.

Other attacks

<i>BinarizationRefinementAttack</i>	For models that preprocess their inputs by binarizing the inputs, this attack can improve adversarials found by other attacks.
<i>PrecomputedImagesAttack</i>	Attacks a model using precomputed adversarial candidates.

1.13 foolbox.adversarial

Provides a class that represents an adversarial example.

```
class foolbox.adversarial.Adversarial (model, criterion, original_image, original_class, distance=<class 'foolbox.distances.MeanSquaredDistance'>, threshold=None, verbose=False)
```

Defines an adversarial that should be found and stores the result.

The *Adversarial* class represents a single adversarial example for a given model, criterion and reference image. It can be passed to an adversarial attack to find the actual adversarial.

Parameters

model [a `Model` instance] The model that should be fooled by the adversarial.

criterion [a `Criterion` instance] The criterion that determines which images are adversarial.

original_image [a `numpy.ndarray`] The original image to which the adversarial image should be as close as possible.

original_class [int] The ground-truth label of the original image.

distance [a `Distance` class] The measure used to quantify similarity between images.

threshold [float or `Distance`] If not `None`, the attack will stop as soon as the adversarial perturbation has a size smaller than this threshold. Can be an instance of the `Distance` class passed to the `distance` argument, or a float assumed to have the same unit as the the given distance. If `None`, the attack will simply minimize the distance as good as possible. Note that the threshold only influences early stopping of the attack; the returned adversarial does not necessarily have smaller perturbation size than this threshold; the `reached_threshold()` method can be used to check if the threshold has been reached.

adversarial_class

The argmax of the model predictions for the best adversarial found so far.

`None` if no adversarial has been found.

backward (*gradient*, *image=None*, *strict=True*)

Interface to `model.backward` for attacks.

Parameters

gradient [`numpy.ndarray`] Gradient of some loss w.r.t. the logits.

image [`numpy.ndarray`] Single input with shape as expected by the model (without the batch dimension).

Returns

gradient [`numpy.ndarray`] The gradient w.r.t the image.

See also:

`gradient()`

batch_predictions (*images*, *greedy=False*, *strict=True*, *return_details=False*)

Interface to `model.batch_predictions` for attacks.

Parameters

images [`numpy.ndarray`] Batch of inputs with shape as expected by the model.

greedy [bool] Whether the first adversarial should be returned.

strict [bool] Controls if the bounds for the pixel values should be checked.

channel_axis (*batch*)

Interface to `model.channel_axis` for attacks.

Parameters

batch [bool] Controls whether the index of the axis for a batch of images (4 dimensions) or a single image (3 dimensions) should be returned.

distance

The distance of the adversarial input to the original input.

gradient (*image=None*, *label=None*, *strict=True*)

Interface to `model.gradient` for attacks.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension). Defaults to the original image.

label [int] Label used to calculate the loss that is differentiated. Defaults to the original label.

strict [bool] Controls if the bounds for the pixel values should be checked.

has_gradient ()

Returns true if `_backward` and `_forward_backward` can be called by an attack, False otherwise.

image

The best adversarial found so far.

normalized_distance (*image*)

Calculates the distance of a given image to the original image.

Parameters

image [*numpy.ndarray*] The image that should be compared to the original image.

Returns

:class:'Distance' The distance between the given image and the original image.

original_class

The class of the original input (ground-truth, not model prediction).

original_image

The original input.

output

The model predictions for the best adversarial found so far.

None if no adversarial has been found.

predictions (*image*, *strict=True*, *return_details=False*)

Interface to `model.predictions` for attacks.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension).

strict [bool] Controls if the bounds for the pixel values should be checked.

predictions_and_gradient (*image=None*, *label=None*, *strict=True*, *return_details=False*)

Interface to `model.predictions_and_gradient` for attacks.

Parameters

image [*numpy.ndarray*] Single input with shape as expected by the model (without the batch dimension). Defaults to the original image.

label [int] Label used to calculate the loss that is differentiated. Defaults to the original label.

strict [bool] Controls if the bounds for the pixel values should be checked.

reached_threshold ()

Returns True if a threshold is given and the currently best adversarial distance is smaller than the threshold.

target_class ()

Interface to `criterion.target_class` for attacks.

1.14 foolbox.utils

`foolbox.utils.softmax` (*logits*)

Transforms predictions into probability values.

Parameters

logits [array_like] The logits predicted by the model.

Returns

‘`numpy.ndarray`’ Probability values corresponding to the logits.

`foolbox.utils.crossentropy` (*label, logits*)

Calculates the cross-entropy.

Parameters

logits [array_like] The logits predicted by the model.

label [int] The label describing the target distribution.

Returns

float The cross-entropy between `softmax(logits)` and `onehot(label)`.

`foolbox.utils.batch_crossentropy` (*label, logits*)

Calculates the cross-entropy for a batch of logits.

Parameters

logits [array_like] The logits predicted by the model for a batch of inputs.

label [int] The label describing the target distribution.

Returns

np.ndarray The cross-entropy between `softmax(logits[i])` and `onehot(label)` for all *i*.

`foolbox.utils.binarize` (*x, values, threshold=None, included_in='upper'*)

Binarizes the values of *x*.

Parameters

values [tuple of two floats] The lower and upper value to which the inputs are mapped.

threshold [float] The threshold; defaults to $(\text{values}[0] + \text{values}[1]) / 2$ if `None`.

included_in [str] Whether the threshold value itself belongs to the lower or upper interval.

`foolbox.utils.imagenet_example` (*shape=(224, 224), data_format='channels_last'*)

Returns an example image and its imagenet class label.

Parameters

shape [list of integers] The shape of the returned image.

data_format [str] “channels_first” or “channels_last”

Returns

image [array_like] The example image.

label [int] The imagenet label associated with the image.

NOTE: This function is deprecated and will be removed in the future.

```
foolbox.utils.samples(dataset='imagenet', index=0, batchsize=1, shape=(224, 224),  
                      data_format='channels_last')
```

Returns a batch of example images and the corresponding labels

Parameters

dataset [string] The data set to load (options: imagenet, mnist, cifar10, cifar100, fashionM-NIST)

index [int] For each data set 20 example images exist. The returned batch contains the images with index [index, index + 1, index + 2, ...]

batchsize [int] Size of batch.

shape [list of integers] The shape of the returned image (only relevant for Imagenet).

data_format [str] "channels_first" or "channels_last"

Returns

images [array_like] The batch of example images

labels [array of int] The labels associated with the images.

```
foolbox.utils.onehot_like(a, index, value=1)
```

Creates an array like a, with all values set to 0 except one.

Parameters

a [array_like] The returned one-hot array will have the same shape and dtype as this array

index [int] The index that should be set to *value*

value [single value compatible with a.dtype] The value to set at the given index

Returns

'numpy.ndarray' One-hot array with the given value at the given location and zeros everywhere else.

CHAPTER 2

Indices and tables

- `genindex`
- `modindex`
- `search`

Bibliography

- [1] Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy, “Explaining and Harnessing Adversarial Examples”, <https://arxiv.org/abs/1412.6572>
- [1] Alexey Kurakin, Ian Goodfellow, Samy Bengio, “Adversarial examples in the physical world”, <https://arxiv.org/abs/1607.02533>
- [1] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, Adrian Vladu, “Towards Deep Learning Models Resistant to Adversarial Attacks”, <https://arxiv.org/abs/1706.06083>
- [1] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, Adrian Vladu, “Towards Deep Learning Models Resistant to Adversarial Attacks”, <https://arxiv.org/abs/1706.06083>
- [1] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, Jianguo Li, “Boosting Adversarial Attacks with Momentum”, <https://arxiv.org/abs/1710.06081>
- [1] <https://arxiv.org/abs/1510.05328>
- [1] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Pascal Frossard, “DeepFool: a simple and accurate method to fool deep neural networks”, <https://arxiv.org/abs/1511.04599>
- [1] Uyeong Jang et al., “Objective Metrics and Gradient Descent Algorithms for Adversarial Examples in Machine Learning”, <https://dl.acm.org/citation.cfm?id=3134635>
- [1] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, Ananthram Swami, “The Limitations of Deep Learning in Adversarial Settings”, <https://arxiv.org/abs/1511.07528>
- [1] Nicholas Carlini, David Wagner: “Towards Evaluating the Robustness of Neural Networks”, <https://arxiv.org/abs/1608.04644>
- [2] https://github.com/carlini/nn_robust_attacks
- [1] Nina Narodytska, Shiva Prasad Kasiviswanathan, “Simple Black-Box Adversarial Perturbations for Deep Networks”, <https://arxiv.org/abs/1612.06299>
- [1] Wieland Brendel (*), Jonas Rauber (*), Matthias Bethge, “Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models”, <https://arxiv.org/abs/1712.04248>
- [1] Logan Engstrom*, Brandon Tran*, Dimitris Tsipras*, Ludwig Schmidt, Aleksander Madry: “A Rotation and a Translation Suffice: Fooling CNNs with Simple Transformations”, <http://arxiv.org/abs/1712.02779>
- [1] L. Schott, J. Rauber, M. Bethge, W. Brendel: “Towards the first adversarially robust neural network model on MNIST”, ICLR (2019) <https://arxiv.org/abs/1805.09190>

f

foolbox.adversarial, 59
foolbox.attacks, 36
foolbox.criteria, 28
foolbox.distances, 35
foolbox.models, 11
foolbox.utils, 62
foolbox.zoo, 33

Symbols

- `__call__()` (*foolbox.attacks.ADefAttack* method), 45
`__call__()` (*foolbox.attacks.AdditiveGaussianNoiseAttack* method), 54
`__call__()` (*foolbox.attacks.AdditiveUniformNoiseAttack* method), 52
`__call__()` (*foolbox.attacks.BinarizationRefinementAttack* method), 56
`__call__()` (*foolbox.attacks.BlendedUniformNoiseAttack* method), 56
`__call__()` (*foolbox.attacks.BoundaryAttack* method), 49
`__call__()` (*foolbox.attacks.CarliniWagnerL2Attack* method), 47
`__call__()` (*foolbox.attacks.ContrastReductionAttack* method), 52
`__call__()` (*foolbox.attacks.DeepFoolAttack* method), 43
`__call__()` (*foolbox.attacks.DeepFoolL2Attack* method), 44
`__call__()` (*foolbox.attacks.DeepFoolLinfinityAttack* method), 44
`__call__()` (*foolbox.attacks.GaussianBlurAttack* method), 51
`__call__()` (*foolbox.attacks.GradientAttack* method), 36
`__call__()` (*foolbox.attacks.GradientSignAttack* method), 36
`__call__()` (*foolbox.attacks.IterativeGradientAttack* method), 46
`__call__()` (*foolbox.attacks.IterativeGradientSignAttack* method), 46
`__call__()` (*foolbox.attacks.L1BasicIterativeAttack* method), 38
`__call__()` (*foolbox.attacks.L2BasicIterativeAttack* method), 39
`__call__()` (*foolbox.attacks.LBFGSAttack* method), 42
`__call__()` (*foolbox.attacks.LinfinityBasicIterativeAttack* method), 37
`__call__()` (*foolbox.attacks.LocalSearchAttack* method), 48
`__call__()` (*foolbox.attacks.MomentumIterativeAttack* method), 41
`__call__()` (*foolbox.attacks.NewtonFoolAttack* method), 43
`__call__()` (*foolbox.attacks.PointwiseAttack* method), 51
`__call__()` (*foolbox.attacks.PrecomputedImagesAttack* method), 57
`__call__()` (*foolbox.attacks.ProjectedGradientDescentAttack* method), 39
`__call__()` (*foolbox.attacks.RandomStartProjectedGradientDescentAttack* method), 40
`__call__()` (*foolbox.attacks.SLSQPAttack* method), 45
`__call__()` (*foolbox.attacks.SaliencyMapAttack* method), 45
`__call__()` (*foolbox.attacks.SaltAndPepperNoiseAttack* method), 56
`__call__()` (*foolbox.attacks.SinglePixelAttack* method), 48
`__call__()` (*foolbox.attacks.SpatialAttack* method), 50
`__class__` (*foolbox.attacks.AdditiveGaussianNoiseAttack* attribute), 54
`__class__` (*foolbox.attacks.AdditiveUniformNoiseAttack* attribute), 52
`__delattr__` (*foolbox.attacks.AdditiveGaussianNoiseAttack* attribute), 54
`__delattr__` (*foolbox.attacks.AdditiveUniformNoiseAttack* attribute), 52
`__dir__()` (*foolbox.attacks.AdditiveGaussianNoiseAttack* method), 54
`__dir__()` (*foolbox.attacks.AdditiveUniformNoiseAttack* method), 53
`__eq__` (*foolbox.attacks.AdditiveGaussianNoiseAttack* attribute), 54
`__eq__` (*foolbox.attacks.AdditiveUniformNoiseAttack* attribute), 54

attribute), 53
 __format__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 54
 __format__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __ge__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 54
 __ge__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __getattr__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 54
 __getattr__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __gt__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 54
 __gt__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __hash__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 54
 __hash__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __init__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __init__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __init__ (foolbox.attacks.ApproximateLBFGSAttack method), 49
 __init__ (foolbox.attacks.LBFGSAttack method), 42
 __init__ (foolbox.attacks.PrecomputedImagesAttack method), 57
 __le__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __le__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __lt__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __lt__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __ne__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __ne__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __new__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __new__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __reduce__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __reduce__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __reduce_ex__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __reduce_ex__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __repr__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __repr__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __setattr__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __setattr__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __sizeof__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __sizeof__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __str__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __str__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53
 __subclasshook__ (foolbox.attacks.AdditiveGaussianNoiseAttack method), 55
 __subclasshook__ (foolbox.attacks.AdditiveUniformNoiseAttack method), 53
 __weakref__ (foolbox.attacks.AdditiveGaussianNoiseAttack attribute), 55
 __weakref__ (foolbox.attacks.AdditiveUniformNoiseAttack attribute), 53

A

AdditiveGaussianNoiseAttack (class in foolbox.attacks), 54
 AdditiveUniformNoiseAttack (class in foolbox.attacks), 52
 ADefAttack (class in foolbox.attacks), 44
 Adversarial (class in foolbox.adversarial), 59
 adversarial_class (foolbox.adversarial.Adversarial attribute), 60
 ApproximateLBFGSAttack (class in foolbox.attacks), 49

B

backward (foolbox.adversarial.Adversarial method), 60

- `backward()` (*foolbox.models.CompositeModel method*), 27
- `backward()` (*foolbox.models.DifferentiableModel method*), 13
- `backward()` (*foolbox.models.KerasModel method*), 19
- `backward()` (*foolbox.models.LasagneModel method*), 22
- `backward()` (*foolbox.models.MXNetGluonModel method*), 24
- `backward()` (*foolbox.models.MXNetModel method*), 23
- `backward()` (*foolbox.models.PyTorchModel method*), 18
- `backward()` (*foolbox.models.TensorFlowEagerModel method*), 16
- `backward()` (*foolbox.models.TensorFlowModel method*), 14
- `backward()` (*foolbox.models.TheanoModel method*), 20
- `BasicIterativeMethod` (*in module foolbox.attacks*), 38
- `batch_crossentropy()` (*in module foolbox.utils*), 62
- `batch_predictions()` (*foolbox.adversarial.Adversarial method*), 60
- `batch_predictions()` (*foolbox.models.CompositeModel method*), 27
- `batch_predictions()` (*foolbox.models.KerasModel method*), 19
- `batch_predictions()` (*foolbox.models.LasagneModel method*), 22
- `batch_predictions()` (*foolbox.models.Model method*), 12
- `batch_predictions()` (*foolbox.models.ModelWrapper method*), 25
- `batch_predictions()` (*foolbox.models.MXNetGluonModel method*), 25
- `batch_predictions()` (*foolbox.models.MXNetModel method*), 23
- `batch_predictions()` (*foolbox.models.PyTorchModel method*), 18
- `batch_predictions()` (*foolbox.models.TensorFlowEagerModel method*), 17
- `batch_predictions()` (*foolbox.models.TensorFlowModel method*), 15
- `batch_predictions()` (*foolbox.models.TheanoModel method*), 20
- `best_other_class()` (*foolbox.attacks.CarliniWagnerL2Attack method*), 47
- `BIM` (*in module foolbox.attacks*), 38
- `BinarizationRefinementAttack` (*class in foolbox.attacks*), 56
- `binarize()` (*in module foolbox.utils*), 62
- `BlendedUniformNoiseAttack` (*class in foolbox.attacks*), 56
- `BoundaryAttack` (*class in foolbox.attacks*), 49
- ## C
- `CarliniWagnerL2Attack` (*class in foolbox.attacks*), 47
- `channel_axis()` (*foolbox.adversarial.Adversarial method*), 60
- `CompositeModel` (*class in foolbox.models*), 26
- `ConfidentMisclassification` (*class in foolbox.criteria*), 30
- `ContrastReductionAttack` (*class in foolbox.attacks*), 52
- `Criterion` (*class in foolbox.criteria*), 29
- `crossentropy()` (*in module foolbox.utils*), 62
- ## D
- `DeepFoolAttack` (*class in foolbox.attacks*), 43
- `DeepFoolL2Attack` (*class in foolbox.attacks*), 44
- `DeepFoolLinfinityAttack` (*class in foolbox.attacks*), 44
- `DifferentiableModel` (*class in foolbox.models*), 13
- `DifferentiableModelWrapper` (*class in foolbox.models*), 26
- `Distance` (*class in foolbox.distances*), 35
- `distance` (*foolbox.adversarial.Adversarial attribute*), 60
- ## F
- `fetch_weights()` (*in module foolbox.zoo*), 34
- `FGSM` (*in module foolbox.attacks*), 37
- `foolbox.adversarial` (*module*), 59
- `foolbox.attacks` (*module*), 36
- `foolbox.criteria` (*module*), 28
- `foolbox.distances` (*module*), 35
- `foolbox.models` (*module*), 11
- `foolbox.utils` (*module*), 62
- `foolbox.zoo` (*module*), 33
- `from_keras()` (*foolbox.models.TensorFlowModel class method*), 15
- ## G
- `GaussianBlurAttack` (*class in foolbox.attacks*), 51
- `get_model()` (*in module foolbox.zoo*), 33
- `gradient()` (*foolbox.adversarial.Adversarial method*), 60
- `gradient()` (*foolbox.models.CompositeModel method*), 27
- `gradient()` (*foolbox.models.DifferentiableModel method*), 13

`gradient()` (*foolbox.models.LasagneModel* method), 22

`gradient()` (*foolbox.models.TensorFlowModel* method), 15

`gradient()` (*foolbox.models.TheanoModel* method), 20

`GradientAttack` (*class in foolbox.attacks*), 36

`GradientSignAttack` (*class in foolbox.attacks*), 36

H

`has_gradient()` (*foolbox.adversarial.Adversarial* method), 61

I

`image` (*foolbox.adversarial.Adversarial* attribute), 61

`imagenet_example()` (*in module foolbox.utils*), 62

`is_adversarial()` (*foolbox.criteria.ConfidentMisclassification* method), 30

`is_adversarial()` (*foolbox.criteria.Criterion* method), 29

`is_adversarial()` (*foolbox.criteria.Misclassification* method), 30

`is_adversarial()` (*foolbox.criteria.OriginalClassProbability* method), 32

`is_adversarial()` (*foolbox.criteria.TargetClass* method), 32

`is_adversarial()` (*foolbox.criteria.TargetClassProbability* method), 33

`is_adversarial()` (*foolbox.criteria.TopKMisclassification* method), 31

`IterativeGradientAttack` (*class in foolbox.attacks*), 46

`IterativeGradientSignAttack` (*class in foolbox.attacks*), 46

K

`KerasModel` (*class in foolbox.models*), 18

L

`L0` (*class in foolbox.distances*), 35

`L1BasicIterativeAttack` (*class in foolbox.attacks*), 38

`L2BasicIterativeAttack` (*class in foolbox.attacks*), 38

`LasagneModel` (*class in foolbox.models*), 21

`LBFGSAttack` (*class in foolbox.attacks*), 42

`Linf` (*in module foolbox.distances*), 36

`Linfinity` (*class in foolbox.distances*), 35

`LinfinityBasicIterativeAttack` (*class in foolbox.attacks*), 37

`LocalSearchAttack` (*class in foolbox.attacks*), 48

`loss_function()` (*foolbox.attacks.CarliniWagnerL2Attack* class method), 47

M

`MAE` (*in module foolbox.distances*), 35

`MeanAbsoluteDistance` (*class in foolbox.distances*), 35

`MeanSquaredDistance` (*class in foolbox.distances*), 35

`Misclassification` (*class in foolbox.criteria*), 29

`Model` (*class in foolbox.models*), 12

`ModelWithEstimatedGradients` (*class in foolbox.models*), 26

`ModelWithoutGradients` (*class in foolbox.models*), 26

`ModelWrapper` (*class in foolbox.models*), 25

`MomentumIterativeAttack` (*class in foolbox.attacks*), 41

`MomentumIterativeMethod` (*in module foolbox.attacks*), 42

`MSE` (*in module foolbox.distances*), 35

`MXNetGluonModel` (*class in foolbox.models*), 24

`MXNetModel` (*class in foolbox.models*), 23

N

`name()` (*foolbox.attacks.AdditiveGaussianNoiseAttack* method), 55

`name()` (*foolbox.attacks.AdditiveUniformNoiseAttack* method), 54

`name()` (*foolbox.attacks.LBFGSAttack* method), 42

`name()` (*foolbox.criteria.ConfidentMisclassification* method), 30

`name()` (*foolbox.criteria.Criterion* method), 29

`name()` (*foolbox.criteria.Misclassification* method), 30

`name()` (*foolbox.criteria.OriginalClassProbability* method), 32

`name()` (*foolbox.criteria.TargetClass* method), 32

`name()` (*foolbox.criteria.TargetClassProbability* method), 33

`name()` (*foolbox.criteria.TopKMisclassification* method), 31

`NewtonFoolAttack` (*class in foolbox.attacks*), 43

`normalized_distance()` (*foolbox.adversarial.Adversarial* method), 61

`num_classes()` (*foolbox.models.CompositeModel* method), 27

`num_classes()` (*foolbox.models.KerasModel* method), 19

`num_classes()` (*foolbox.models.LasagneModel* method), 22

`num_classes()` (*foolbox.models.Model* method), 13

- `num_classes()` (*foolbox.models.ModelWrapper method*), 26
`num_classes()` (*foolbox.models.MXNetGluonModel method*), 25
`num_classes()` (*foolbox.models.MXNetModel method*), 24
`num_classes()` (*foolbox.models.PyTorchModel method*), 18
`num_classes()` (*foolbox.models.TensorFlowEagerModel method*), 17
`num_classes()` (*foolbox.models.TensorFlowModel method*), 16
`num_classes()` (*foolbox.models.TheanoModel method*), 21
- O**
- `onehot_like()` (*in module foolbox.utils*), 63
`original_class` (*foolbox.adversarial.Adversarial attribute*), 61
`original_image` (*foolbox.adversarial.Adversarial attribute*), 61
`OriginalClassProbability` (*class in foolbox.criteria*), 32
`output` (*foolbox.adversarial.Adversarial attribute*), 61
- P**
- `PointwiseAttack` (*class in foolbox.attacks*), 51
`PrecomputedImagesAttack` (*class in foolbox.attacks*), 57
`predictions()` (*foolbox.adversarial.Adversarial method*), 61
`predictions()` (*foolbox.models.Model method*), 13
`predictions()` (*foolbox.models.ModelWrapper method*), 26
`predictions_and_gradient()` (*foolbox.adversarial.Adversarial method*), 61
`predictions_and_gradient()` (*foolbox.models.CompositeModel method*), 27
`predictions_and_gradient()` (*foolbox.models.DifferentiableModel method*), 14
`predictions_and_gradient()` (*foolbox.models.KerasModel method*), 19
`predictions_and_gradient()` (*foolbox.models.LasagneModel method*), 22
`predictions_and_gradient()` (*foolbox.models.MXNetGluonModel method*), 25
`predictions_and_gradient()` (*foolbox.models.MXNetModel method*), 24
`predictions_and_gradient()` (*foolbox.models.PyTorchModel method*), 18
`predictions_and_gradient()` (*foolbox.models.TensorFlowEagerModel method*), 17
`predictions_and_gradient()` (*foolbox.models.TensorFlowModel method*), 16
`predictions_and_gradient()` (*foolbox.models.TheanoModel method*), 21
`ProjectedGradientDescent` (*in module foolbox.attacks*), 40
`ProjectedGradientDescentAttack` (*class in foolbox.attacks*), 39
`PyTorchModel` (*class in foolbox.models*), 17
- R**
- `RandomPGD` (*in module foolbox.attacks*), 41
`RandomProjectedGradientDescent` (*in module foolbox.attacks*), 41
`RandomStartProjectedGradientDescentAttack` (*class in foolbox.attacks*), 40
`reached_threshold()` (*foolbox.adversarial.Adversarial method*), 61
- S**
- `SaliencyMapAttack` (*class in foolbox.attacks*), 45
`SaltAndPepperNoiseAttack` (*class in foolbox.attacks*), 55
`samples()` (*in module foolbox.utils*), 62
`SinglePixelAttack` (*class in foolbox.attacks*), 48
`SLSQPAttack` (*class in foolbox.attacks*), 45
`softmax()` (*in module foolbox.utils*), 62
`SpatialAttack` (*class in foolbox.attacks*), 50
- T**
- `target_class()` (*foolbox.adversarial.Adversarial method*), 61
`TargetClass` (*class in foolbox.criteria*), 31
`TargetClassProbability` (*class in foolbox.criteria*), 33
`TensorFlowEagerModel` (*class in foolbox.models*), 16
`TensorFlowModel` (*class in foolbox.models*), 14
`TheanoModel` (*class in foolbox.models*), 20
`TopKMisclassification` (*class in foolbox.criteria*), 31