
Foolbox Documentation

Release 0.9.5

Jonas Rauber & Wieland Brendel

Aug 17, 2017

1	Installation	3
1.1	Stable release	3
1.2	Development version	3
1.3	Contributing to Foolbox	3
2	Tutorial	5
2.1	Creating a model	5
2.2	Specifying the criterion	5
2.3	Running the attack	6
2.4	Visualizing the adversarial examples	6
3	Examples	7
3.1	Creating a model	7
3.2	Applying an attack	8
4	Advanced	11
4.1	Implicit	11
4.2	Explicit	11
5	Development	13
5.1	New Adversarial Attacks	13
6	foolbox.models	15
6.1	Models	15
6.2	Wrappers	15
6.3	Detailed description	16
7	foolbox.criteria	21
7.1	Criteria	21
7.2	Examples	21
7.3	Detailed description	22
8	foolbox.distances	25
8.1	Distances	25
8.2	Aliases	25
8.3	Base class	25
8.4	Detailed description	25

9	foolbox.attacks	27
9.1	Gradient-based attacks	27
9.2	Black-box attacks	28
9.3	Other attacks	29
10	foolbox.adversarial	31
11	foolbox.utils	33
12	Indices and tables	35
	Bibliography	37
	Python Module Index	39

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.

Foolbox is work in progress and any input is welcome. The source code can be found on [GitHub](#).

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 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*.

1.2 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
```

1.3 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.

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

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))
```

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 TopKMisclassification

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

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=image, label=label)
```

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)
```

Here you can find a collection of examples how Foolbox models can be created using different deep learning frameworks.

3.1 Creating a model

3.1.1 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)
model = foolbox.models.KerasModel(kmodel, bounds=(0, 255),
    ↪preprocessing=preprocessing)

image, _ = foolbox.utils.imagenet_example()
print(np.argmax(model.predictions(image)))
```

3.1.2 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(images, 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()
```

3.2 Applying an attack

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

3.2.1 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)
```


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.

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=image, label=label, unpack=False)
```

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

```
adversarial_image = adversarial.image
```

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)
```


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

5.1 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 and implement the `_apply()` method. The `_apply()` method 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.

 foolbox.models

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

6.1 Models

<i>Model</i>	Base class to provide attacks with a unified interface to models.
<i>DifferentiableModel</i>	Base class for differentiable models that provide gradients.
<i>TensorFlowModel</i>	Creates a <i>Model</i> instance from existing <i>TensorFlow</i> tensors.
<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.

6.2 Wrappers

<i>ModelWrapper</i>	Base class for models that wrap other models.
<i>GradientLess</i>	Turns a model into a model without gradients.
<i>CompositeModel</i>	Combines predictions of a (black-box) model with the gradient of a (substitute) model.

6.3 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 images with shape (batch size, height, width, channels).

Returns *numpy.ndarray*

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

See also:

`prediction()`

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*

Image with shape (height, width, channels).

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.

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*

Image with shape (height, width, channels).

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*

Image with shape (height, width, channels).

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.

class foolbox.models.**PyTorchModel** (*model, bounds, num_classes, channel_axis=1, cuda=True, 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.

cuda : bool

A boolean specifying whether the model uses CUDA.

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.

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.

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

Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

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

Elementwise preprocessing of input; we first subtract the first element of preprocessing from the input and then divide the input by the second element.

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.

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

Turns a model into a model without gradients.

Provides classes that define what is adversarial.

7.1 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>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.

7.2 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
```

7.3 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.

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.

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.

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 $0 \leq p \leq 1$.

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 $0 \leq p \leq 1$.

 foolbox.distances

Provides classes to measure the distance between images.

8.1 Distances

<i>MeanSquaredDistance</i>	Calculates the mean squared error between two images.
<i>MeanAbsoluteDistance</i>	Calculates the mean absolute error between two images.

8.2 Aliases

<i>MSE</i>	alias of <i>MeanSquaredDistance</i>
------------	-------------------------------------

8.3 Base class

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

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

8.4 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.

foolbox.distances.**MSE**

alias of *MeanSquaredDistance*

9.1 Gradient-based attacks

```
class foolbox.attacks.GradientSignAttack (model=None, crite-  

riterion=<foolbox.criteria.Misclassification ob-  

ject>)
```

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

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

```
class foolbox.attacks.IterativeGradientSignAttack (model=None, crite-  

riterion=<foolbox.criteria.Misclassification  

object>)
```

Like GradientSignAttack but with several steps for each epsilon.

```
class foolbox.attacks.GradientAttack (model=None, criterion=<foolbox.criteria.Misclassification  

object>)
```

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.

```
class foolbox.attacks.IterativeGradientAttack (model=None, crite-  

riterion=<foolbox.criteria.Misclassification  

object>)
```

Like GradientAttack but with several steps for each epsilon.

```
foolbox.attacks.FGSM  

  alias of GradientSignAttack
```

```
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 [R11] to support other targeted criteria and other distance measures.

References

[R11]

class foolbox.attacks.**DeepFoolAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

Simple and accurate adversarial attack.

Implements DeepFool introduced in [R33].

References

[R33]

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

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

class foolbox.attacks.**SaliencyMapAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

Implements the Saliency Map Attack.

The attack was introduced in [R55].

References

[R55]

9.2 Black-box attacks

class foolbox.attacks.**GaussianBlurAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

Blurs the image until it is misclassified.

class foolbox.attacks.**ContrastReductionAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

Reduces the contrast of the image until it is misclassified.

class foolbox.attacks.**SinglePixelAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

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

class foolbox.attacks.**LocalSearchAttack** (*model=None, criterion=<foolbox.criteria.Misclassification object>*)

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

This implementation is based on the algorithm in [R99].

References

[R99]

class foolbox.attacks.**AdditiveUniformNoiseAttack** (*model=None*, *criterion=<foolbox.criteria.Misclassification object>*)

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

class foolbox.attacks.**AdditiveGaussianNoiseAttack** (*model=None*, *criterion=<foolbox.criteria.Misclassification object>*)

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

class foolbox.attacks.**SaltAndPepperNoiseAttack** (*model=None*, *criterion=<foolbox.criteria.Misclassification object>*)

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

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

Same as LBFGSAttack with approximate_gradient set to True.

9.3 Other attacks

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.

Gradient-based attacks

<i>GradientSignAttack</i>	Adds the sign of the gradient to the image, gradually increasing the magnitude until the image is misclassified.
<i>IterativeGradientSignAttack</i>	Like GradientSignAttack but with several steps for each epsilon.
<i>GradientAttack</i>	Perturbs the image with the gradient of the loss w.r.t.
<i>IterativeGradientAttack</i>	Like GradientAttack but with several steps for each epsilon.
<i>FGSM</i>	alias of <i>GradientSignAttack</i>
<i>LBFGSAttack</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 one-hot encoded target class.
<i>DeepFoolAttack</i>	Simple and accurate adversarial attack.
<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.

Black-box attacks

<i>GaussianBlurAttack</i>	Blurs the image until it is misclassified.
<i>ContrastReductionAttack</i>	Reduces the contrast of the image until it is misclassified.
<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>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>ApproximateLBFSGSAttack</i>	Same as LBFSGSAttack with approximate_gradient set to True.

Other attacks

<i>PrecomputedImagesAttack</i>	Attacks a model using precomputed adversarial candidates.
--------------------------------	---

Provides a class that represents an adversarial example.

```
class foolbox.adversarial.Adversarial (model, criterion, original_image, original_class,
                                       distance=<class 'foolbox.distances.MeanSquaredDistance'>,
                                       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.

batch_predictions (*images*, *increasing*=`False`, *strict*=`True`)

Interface to `model.batch_predictions` for attacks.

Parameters *images* : `numpy.ndarray`

Batch of images with shape (batch size, height, width, channels).

increasing : `bool`

Whether the images can be assumed to have increasing distance from the original image.

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.

gradient (*image=None, label=None, strict=True*)

Interface to `model.gradient` for attacks.

Parameters **image** : *numpy.ndarray*

Image with shape (height, width, channels). 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.

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 *Distance*

The distance between the given image and the original image.

predictions (*image, strict=True*)

Interface to `model.predictions` for attacks.

Parameters **image** : *numpy.ndarray*

Image with shape (height, width, channels).

strict : bool

Controls if the bounds for the pixel values should be checked.

predictions_and_gradient (*image=None, label=None, strict=True*)

Interface to `model.predictions_and_gradient` for attacks.

Parameters **image** : *numpy.ndarray*

Image with shape (height, width, channels). 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.

target_class ()

Interface to `criterion.target_class` for attacks.

CHAPTER 11

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).

CHAPTER 12

Indices and tables

- `genindex`
- `modindex`
- `search`

Bibliography

- [R11] <https://arxiv.org/abs/1510.05328>
- [R33] DeepFool: a simple and accurate method to fool deep neural networks Seyed-Mohsen Moosavi-Dezfooli, Al-hussein Fawzi, Pascal Frossard <https://arxiv.org/abs/1511.04599>
- [R55] The Limitations of Deep Learning in Adversarial Settings Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, Ananthram Swami <https://arxiv.org/abs/1511.07528>
- [R99] Simple Black-Box Adversarial Perturbations for Deep Networks Nina Narodytska, Shiva Prasad Kasiviswanathan <https://arxiv.org/pdf/1612.06299.pdf>

f

foolbox.adversarial, 31
foolbox.attacks, 27
foolbox.criteria, 21
foolbox.distances, 25
foolbox.models, 15
foolbox.utils, 33

A

AdditiveGaussianNoiseAttack (class in foolbox.attacks), 29
AdditiveUniformNoiseAttack (class in foolbox.attacks), 29
Adversarial (class in foolbox.adversarial), 31
ApproximateLBFGSAttack (class in foolbox.attacks), 29

B

batch_predictions() (foolbox.adversarial.Adversarial method), 31
batch_predictions() (foolbox.models.Model method), 16

C

channel_axis() (foolbox.adversarial.Adversarial method), 32
ContrastReductionAttack (class in foolbox.attacks), 28
Criterion (class in foolbox.criteria), 22
crossentropy() (in module foolbox.utils), 33

D

DeepFoolAttack (class in foolbox.attacks), 28
DifferentiableModel (class in foolbox.models), 16
Distance (class in foolbox.distances), 25

F

FGSM (in module foolbox.attacks), 27
foolbox.adversarial (module), 31
foolbox.attacks (module), 27
foolbox.criteria (module), 21
foolbox.distances (module), 25
foolbox.models (module), 15
foolbox.utils (module), 33

G

GaussianBlurAttack (class in foolbox.attacks), 28
gradient() (foolbox.adversarial.Adversarial method), 32
gradient() (foolbox.models.DifferentiableModel method), 17

GradientAttack (class in foolbox.attacks), 27
GradientLess (class in foolbox.models), 19
GradientSignAttack (class in foolbox.attacks), 27

H

has_gradient() (foolbox.adversarial.Adversarial method), 32

I

is_adversarial() (foolbox.criteria.Criterion method), 22
IterativeGradientAttack (class in foolbox.attacks), 27
IterativeGradientSignAttack (class in foolbox.attacks), 27

K

KerasModel (class in foolbox.models), 18

L

LasagneModel (class in foolbox.models), 19
LBFGSAttack (class in foolbox.attacks), 27
LocalSearchAttack (class in foolbox.attacks), 28

M

MeanAbsoluteDistance (class in foolbox.distances), 26
MeanSquaredDistance (class in foolbox.distances), 25
Misclassification (class in foolbox.criteria), 22
Model (class in foolbox.models), 16
ModelWrapper (class in foolbox.models), 19
MSE (in module foolbox.distances), 26

N

name() (foolbox.criteria.Criterion method), 22
normalized_distance() (foolbox.adversarial.Adversarial method), 32
num_classes() (foolbox.models.Model method), 16

O

OriginalClassProbability (class in foolbox.criteria), 23

P

PrecomputedImagesAttack (class in foolbox.attacks), 29
predictions() (foolbox.adversarial.Adversarial method), 32
predictions() (foolbox.models.Model method), 16
predictions_and_gradient() (foolbox.adversarial.Adversarial method), 32
predictions_and_gradient() (foolbox.models.DifferentiableModel method), 17
PyTorchModel (class in foolbox.models), 18

S

SaliencyMapAttack (class in foolbox.attacks), 28
SaltAndPepperNoiseAttack (class in foolbox.attacks), 29
SinglePixelAttack (class in foolbox.attacks), 28
SLSQPAttack (class in foolbox.attacks), 28
softmax() (in module foolbox.utils), 33

T

target_class() (foolbox.adversarial.Adversarial method), 32
TargetClass (class in foolbox.criteria), 23
TargetClassProbability (class in foolbox.criteria), 23
TensorFlowModel (class in foolbox.models), 17
TheanoModel (class in foolbox.models), 18
TopKMisclassification (class in foolbox.criteria), 23