

Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



ML Security with the Adversarial Robustness Toolbox

Part 1 — Attacking Machine Learning Models



Kedion · [Follow](#)

17 min read · Apr 26, 2022

Listen

Share

More

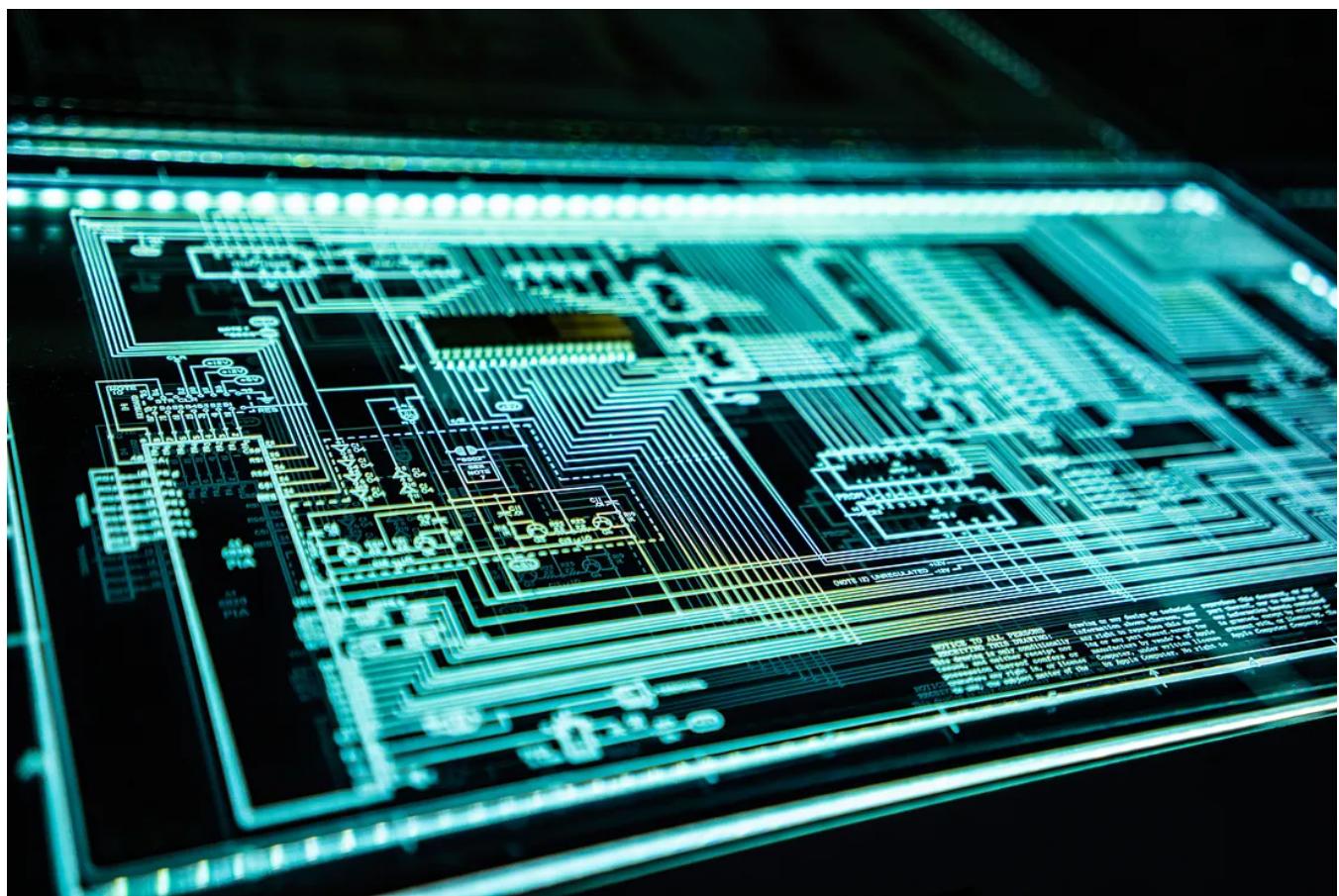


Photo by [Adi Goldstein](#)

Written by [Tigran Avetisyan](#).

Developing machine learning models and putting them into production can be very challenging. However, successfully deploying an ML pipeline is just part of the story – you also need to think about keeping it secure.

Machine learning models are used in a wide range of areas, like finance, medicine, or surveillance, and can be accurate at detecting fraud or filtering out faulty products. In applications like detecting fraud, scammers might be interested in fooling machine learning systems so that they miss scam emails or phishing links.

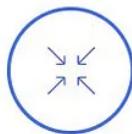
In this series of tutorials, we are going to have a look at the [Adversarial Robustness Toolbox](#) and figure out how it can help you with securing your machine learning pipelines. In PART 1, we will focus on adversarial attacks, and we will be using the MNIST digits dataset along with the TensorFlow/Keras ML framework.

Let's get started!

What is the Adversarial Robustness Toolbox?

The Adversarial Robustness Toolbox, or ART, is a Python framework for machine learning security. ART contains attack and defense tools that can help teams better understand adversarial attacks and develop protection measures based on experimentation.

Features



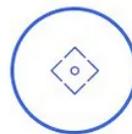
Extended Support

ART supports all popular machine learning frameworks (TensorFlow, Keras, PyTorch, MXNet, scikit-learn, XGBoost, LightGBM, CatBoost, GPy, etc.), all data types (images, tables, audio, video, etc.) and machine learning tasks (classification, object detection, generation, certification, etc.).



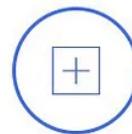
39 Attack Modules

On a high level, ART supports 4 attack modules: Evasion, Poisoning, Extraction, and Inference. Detailed information about the supported attack modules can be found [here](#).



29 Defense Modules

On a high level, ART supports 5 attack modules: Preprocessor, Postprocessor, Trainer, Transformer, and Detector. Detailed information about the supported defense modules can be found [here](#).



Estimators and Metrics

ART supports 3 robustness metrics, 1 certification and 1 verification metric. It also supports multiple estimators and details about the same can be found [here](#).

<https://adversarial-robustness-toolbox.org/>

ART has 39 attack modules, 29 defense modules, and supports a wide range of machine learning frameworks, including scikit-learn, PyTorch, TensorFlow, and Keras. ART also supports several machine learning tasks (including classification, regression, and generation) and works with all data types (audio, video, images, or tables).

The Adversarial Robustness Toolbox was originally developed and published by IBM. In July 2020, IBM donated ART to the Linux Foundation AI (LFAI). Since then, LFAI has maintained and developed updates for the toolkit.

Attack Types in the Adversarial Robustness Toolbox

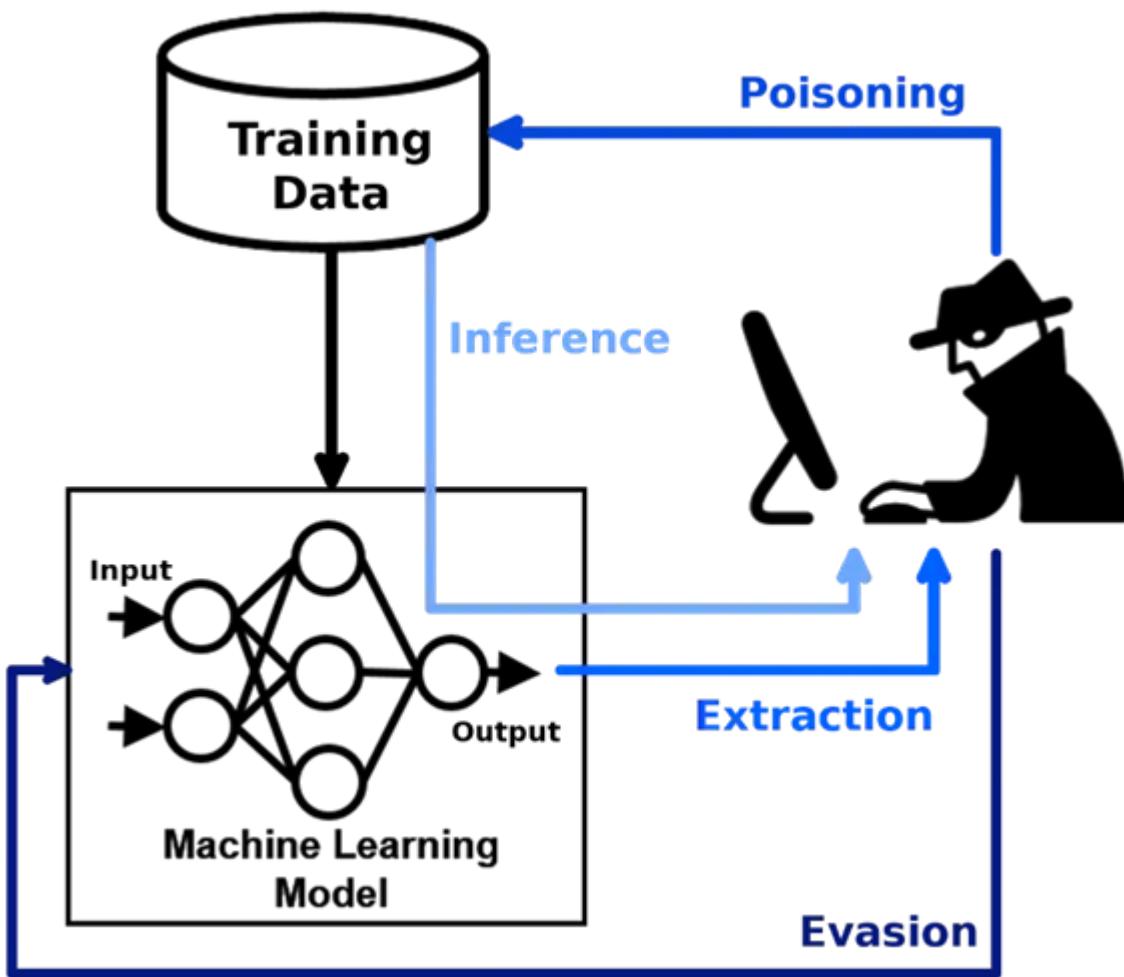
Because PART 1 in this series focuses on attacks, let's take a deeper look at the attack types supported by ART.

At a high level, there are 4 types of adversarial attacks implemented in ART:

- **Evasion.** Evasion attacks typically work by perturbing input data to cause a trained model to misclassify it. Evasion is done after training and during inference, i.e. when models are already deployed in production. Adversaries perform evasion attacks to avoid detection by AI systems. As an example, adversaries might run an evasion attack to cause the victim model to miss phishing emails. Evasion attacks might require access to the victim model.
- **Extraction.** Extraction is an attack where an adversary attempts to build a model that is similar or identical to a victim model. In simple words, extraction is the attempt of copying or stealing a machine learning model. Extraction attacks typically require access to the original model, as well as to data that is similar or identical to the data originally used to train the victim model.
- **Inference.** Inference attacks generally aim at reconstructing a part or the entirety of the dataset that was used to train the victim model. Adversaries can use inference attacks to reconstruct entire training samples, separate features, or determine if a sample has been used to train the victim model. Inference attacks typically require access to the victim model. In some cases, attackers might also need to have access to some portion of the data used to train the model.
- **Poisoning.** Poisoning attacks aim to perturb training data to corrupt the victim model during training. Poisoned data contains features (called a backdoor) that trigger the desired output in a trained model. Essentially, the perturbed features

cause the model to overfit to them. As a very simple example (which we'll have a look at in code below), an attacker could poison the digits in the MNIST dataset so that the victim model classifies all digits as 9s. Poisoning attacks require access to the training data of a model before the actual training occurs.

ART's documentation supplies this neat graph that shows how the attacks work at a high level:



<https://adversarial-robustness-toolbox.readthedocs.io/en/latest/index.html>

For the vast majority of implemented attacks, ART supplies links to the research papers that provide more detail on a given attack. So if you want to learn more about a specific attack, look for paper links in ART's documentation.

Below, we will take a look at each of the attack types supported by ART. We will implement one attack per type, but what you'll see below should transfer to other attacks of the same type as well.

How do attacks on machine learning pipelines happen?

As we pointed out above, adversaries typically need some form of access to your machine learning model or its training data to perform an attack. But assuming that your model is hosted in an environment that an adversary can't reach, how do attacks on ML pipelines even happen? In more practical terms, how can an adversary gain access to your model and training data?

Here are just some of the cases in which adversaries can obtain access to your pipeline and data:

- **You are using a dataset from an unverified source.** An adversary could poison data and publish it somewhere for unsuspecting teams to use. Datasets from untrustworthy sources or datasets that are not verified are at higher risk of being compromised.
- **You are using a model from an unverified source.** Similar to training data, adversaries can publish models for victims to use. They can later use their knowledge of the model to perform attacks on it. Models from suspicious sources might also be pre-trained on poisoned data.
- **You are using a model that is available publicly.** Suppose that you are using an NLP model from Hugging Face. Hugging Face models are available to anyone. An adversary could analyze a publicly available model and use their knowledge to attack other similar models used by other teams. In fact, adversaries don't need to have the exact same model as you do — knowledge of a model that performs the same task as yours can be enough for an attack.
- **An adversary has access to your ML pipeline.** If a malicious insider gains access to your pipeline — by leveraging their position inside the organization, for example — they can get in-depth knowledge about its architecture, weights, and training data. They might also know about the measures that you use to protect your model, which makes malicious insiders arguably the biggest threat to ML pipeline security.

This isn't an exhaustive list, but it should give you an idea of how your model or training data could become compromised.

Prerequisites for Using ART

To be able to follow along with this tutorial, install ART by using this command.

```
pip install adversarial-robustness-toolbox
```

If you are using conda, use this command instead:

```
conda install -c conda-forge adversarial-robustness-toolbox
```

We used ART version 1.10.0, but newer versions should work fine as well.

If necessary, you can install ART just for the specific ML/DL framework that you will be using. You can learn more about this [here](#).

Aside from ART, you will need whichever ML/DL framework you want to attack. We are going to be using TensorFlow, which you can install by using this command:

```
pip install tensorflow
```

Or if you are using conda, this command:

```
conda install -c conda-forge tensorflow
```

We are also going to be using NumPy and Matplotlib, so make sure that you have these libraries as well.

Evasion Attacks in ART

Let's start with evasion attacks in ART. As mentioned earlier, an evasion attack is when the attacker perturbs input at inference time to cause the model to misclassify it.

As of ART version 1.10.0, most of the attacks supported by the framework were evasion attacks.

We will be using the Fast Gradient Method to generate adversarial samples from our test set. You can read more about the Fast Gradient Method in [this paper](#).

Note that you can find the full code for this tutorial [here](#). You can find more examples from the authors of ART [here](#) and [here](#).

Importing dependencies

As usual, we start by importing dependencies:

```
1 # Importing dependencies
2 import tensorflow as tf
3 from tensorflow.keras.layers import Conv2D, MaxPool2D, Dense, Flatten
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import matplotlib
7 from art.estimators.classification import KerasClassifier
8 from art.attacks.evasion import FastGradientMethod
9 from art.utils import load_dataset
10
11 # Disabling eager execution from TF 2
12 tf.compat.v1.disable_eager_execution()
13
14 # Increasing Matplotlib font size
15 matplotlib.rcParams.update({"font.size": 14})
```

dependencies.py hosted with ❤ by GitHub

[view raw](#)

Note that we are disabling eager execution in TensorFlow for this guide on line 12. This is because the wrapper class `art.estimators.classification.KerasClassifier` doesn't fully support TF 2.

Loading data

To load the MNIST dataset, we are going to use ART's function `art.utils.load_dataset`. This function returns tuples for the train and test sets, as well as the minimum and maximum feature values in the dataset.

```
1 # Loading the data
2 (train_images, train_labels), (test_images, test_labels), min, max = load_dataset(name="mnist")
```

The images are already normalized to the `[0, 1]` range, while the labels are one-hot encoded. We don't need to do any preprocessing on the dataset.

Training a TensorFlow Keras model

Now, let's create a simple TensorFlow Keras model:

```
1 # Function for creating model
2 def create_model():
3     # Defining the model
4     model = tf.keras.models.Sequential([
5         Conv2D(filters=32, kernel_size=3, activation="relu", input_shape=(28, 28, 1)),
6         MaxPool2D(pool_size=2),
7         Conv2D(filters=64, kernel_size=3, activation="relu"),
8         MaxPool2D(pool_size=2),
9         Flatten(),
10        Dense(units=10, activation="softmax")
11    ])
12
13    # Compiling the model
14    model.compile(
15        optimizer="adam",
16        loss="categorical_crossentropy",
17        metrics=["accuracy"]
18    )
19
20    # Returning the model
21    return model
```

create_model.py hosted with ❤️ by GitHub

[view raw](#)

And train it:

```
1 # Instantiating the model
2 model = create_model()
3
4 # Training the model
5 model.fit(
6     x=train_images,
7     y=train_labels,
8     epochs=10,
9     batch_size=256)
```

train_model.py hosted with ❤ by GitHub

[view raw](#)

```
1 Train on 60000 samples
2 Epoch 1/10
3 60000/60000 [=====] - 4s 67us/sample - loss: 0.3924 - accuracy: 0.8961
4 Epoch 2/10
5 60000/60000 [=====] - 2s 33us/sample - loss: 0.0934 - accuracy: 0.9721
6 Epoch 3/10
7 60000/60000 [=====] - 2s 33us/sample - loss: 0.0665 - accuracy: 0.9804
8 Epoch 4/10
9 60000/60000 [=====] - 2s 32us/sample - loss: 0.0530 - accuracy: 0.9838
10 Epoch 5/10
11 60000/60000 [=====] - 2s 32us/sample - loss: 0.0467 - accuracy: 0.9858
12 Epoch 6/10
13 60000/60000 [=====] - 2s 32us/sample - loss: 0.0404 - accuracy: 0.9875
14 Epoch 7/10
15 60000/60000 [=====] - 2s 32us/sample - loss: 0.0363 - accuracy: 0.9889
16 Epoch 8/10
17 60000/60000 [=====] - 2s 31us/sample - loss: 0.0322 - accuracy: 0.9903
18 Epoch 9/10
19 60000/60000 [=====] - 2s 31us/sample - loss: 0.0288 - accuracy: 0.9912
20 Epoch 10/10
21 60000/60000 [=====] - 2s 32us/sample - loss: 0.0274 - accuracy: 0.9918
```

If training takes too much time for you, you can try to simplify the model. And if you encounter out-of-memory (OOM) issues, reducing the batch size should be able to help.

Defining an evasion attack on our model

As the next step, let's define an evasion attack for our model.

To be able to run attacks on our model, we must wrap it in the `art.estimators.classification.KerasClassifier` class. Here's how this is done:

```

1 # Creating a classifier by wrapping our TF model in ART's KerasClassifier class
2 classifier = KerasClassifier(
3     model=model,
4     clip_values=(min, max)
5 )

```

[wrap_model.py](#) hosted with ❤ by GitHub

[view raw](#)

The argument `model=model` indicates our model, while `clip_values=(min, max)` specifies the minimum and maximum values allowed for the features. We are using the values provided by the `art.utils.load_dataset` function.

Instead of `KerasClassifier`, you can also try `TensorFlowV2Classifier`, but note that its usage is different.

We then define the attack by using ART's `FastGradientMethod` class:

```

1 # Defining an attack using the fast gradient method
2 attack_fgsm = FastGradientMethod(
3     estimator=classifier,
4     eps=0.3
5 )

```

[fast_gradient_attack.py](#) hosted with ❤ by GitHub

[view raw](#)

`estimator=classifier` shows that the attack will apply to our classifier, while the argument `eps=0.3` essentially defines how strong the attack will be.

We can then generate adversarial samples by calling the attack object's method `generate`, passing to it the target images that we want to perturb:

```

1 # Generating adversarial images from test images
2 test_images_adv = attack_fgsm.generate(x=test_images)

```

[generate_adversarial.py](#) hosted with ❤ by GitHub

[view raw](#)

Evaluating the effectiveness of the attack

Let's take a look at one adversarial sample:

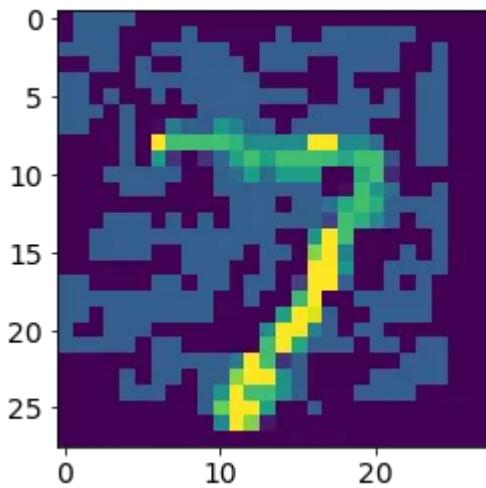
```

1 # Viewing an adversarial image
2 plt.imshow(X=test_images_adv[0])

```

[show_adversarial.py](#) hosted with ❤ by GitHub

[view raw](#)



The fast gradient method applied noise to the clean test images. We can see the effect of the attack by comparing the performance of our model on the clean and adversarial sets:

```

1 # Evaluating the model on clean images
2 score_clean = model.evaluate(
3     x=test_images,
4     y=test_labels
5 )
6
7 # Evaluating the model on adversarial images
8 score_adv = model.evaluate(
9     x=test_images_adv,
10    y=test_labels
11 )
12
13 # Comparing test losses
14 print(f"Clean test set loss: {score_clean[0]:.2f} "
15       f"vs adversarial set test loss: {score_adv[0]:.2f}")
16
17 # Comparing test accuracies
18 print(f"Clean test set accuracy: {score_clean[1]:.2f} "
19       f"vs adversarial test set accuracy: {score_adv[1]:.2f}")

```

[evaluate_evasion.py](#) hosted with ❤ by GitHub

[view raw](#)

Clean test set loss: 0.04 vs adversarial set test loss: 5.95
Clean test set accuracy: 0.99 vs adversarial test set accuracy: 0.07

The attack has affected the model considerably, making it completely unusable. However, because the perturbations in the image are very evident, it would be very

easy for the victim to figure out that they are being attacked.

Lowering the `eps` value can make the tampering less visible, but the impact on the model will be lessened as well. We've tried 10 different values for `eps` to explore their effect on the visual appearance of adversarial samples and on the model's performance.

```
1 # Setting the number of rows and columns for the figure
2 nrows, ncols = 2, 5
3
4 # Generating subplots
5 fig, axes = plt.subplots(
6     nrows=nrows,
7     ncols=ncols,
8     figsize=(20, 10)
9 )
10
11 # Defining a range of eps values to try
12 eps_to_try = [0.01, 0.025, 0.05, 0.075, 0.1, 0.125, 0.15, 0.175, 0.2, 0.25]
13
14 # Defining a counting variable to traverse eps_to_try
15 counter = 0
16
17 # Iterating over rows and cols
18 for i in range(nrows):
19     for j in range(ncols):
20         # Creating an attack object for the current value of eps
21         attack_fgsm = FastGradientMethod(
22             estimator=classifier,
23             eps=eps_to_try[counter]
24         )
25
26         # Generating adversarial images
27         test_images_adv = attack_fgsm.generate(x=test_images)
28
29         # Showing the first adversarial image
30         axes[i, j].imshow(X=test_images_adv[0])
31
32         # Disabling x and y ticks
33         axes[i, j].set_xticks(ticks=[])
34         axes[i, j].set_yticks(ticks=[])
35
36         # Evaluating model performance on adversarial samples and retrieving test accuracy
37         test_score = classifier._model.evaluate(
38             x=test_images_adv,
39             y=test_labels
40             )[1]
41
42         # Getting prediction for the image that we displayed
43         prediction = np.argmax(model.predict(
44             x=np.expand_dims(a=test_images_adv[0],
45             axis=0)
46             ))
47
48         # Showing the current eps value, test accuracy, and prediction
```

```

49         axes[i, j].set_title(
50             label=f"Eps value: {eps_to_try[counter]}\n"
51             f"Test accuracy: {test_score * 100:.2f}%\n"
52             f"Prediction: {prediction}"
53         )
54
55     # Incrementing counter
56     counter += 1
57
58 # Showing the plot
59 plt.show()

```

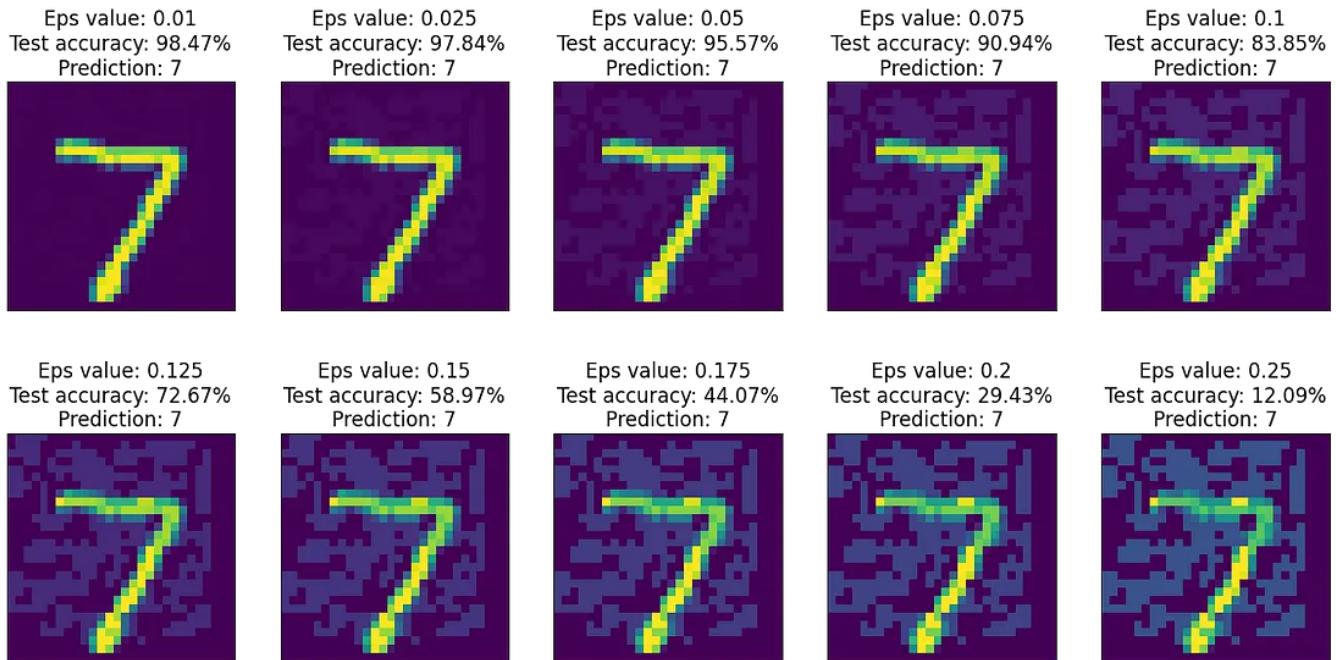
Comments are now hosted with  by GitHub

View raw

In the code block above, we create a figure along with subplots (lines 5 to 9) and define `eps` values to try (line 12). We then iterate over our subplots (lines 18 and 19), generate adversarial samples for each value from `eps_to_try` (lines 21 to 27), plot a sample adversarial image for the current `eps` value (line 30), evaluate model accuracy on the adversarial set (lines 37 to 40), and get a prediction for the sample adversarial image that is being displayed (lines 43 to 46).

Each `eps` value, the test accuracy for it, and the prediction for the sample adversarial image are displayed above each image (lines 49 to 53).

The resulting plot is as follows:



We can see that higher `eps` values produce more visible noise in the image and have a bigger impact on the model's performance. And while the model predicted the

correct label for this particular image 9 times out of 10, the overall performance of the model on the entire adversarial test set was much worse.

Extraction Attacks in ART

As the second step, let's show extraction attacks in ART. Extraction attacks, as a reminder, aim to copy or steal a victim model.

Let's use the class `art.attacks.extraction.CopyCatCNN` to perform the attack. You can learn more about this attack method in [this paper](#).

Training a victim model

For this attack, let's separate our training dataset into two subsets — one with 50,000 samples and the other with 10,000.

```
1 # Importing CopycatCNN
2 from art.attacks.extraction import CopycatCNN
3
4 # Setting aside a subset of the source dataset for the original model
5 train_images_original = train_images[:50000]
6 train_labels_original = train_labels[:50000]
7
8 # Using the rest of the source dataset for the stolen model
9 train_images_stolen = train_images[50000:]
10 train_labels_stolen = train_labels[50000:]
```

define_subsets.py hosted with ❤ by GitHub

[view raw](#)

The subset with 50,000 samples will be used to train the original model, while the subset with 10,000 samples will be used to steal the original model. Basically, we are simulating a situation where an adversary has a dataset that is similar to the original dataset.

Let's train our original model on its dataset:

```

1 # Training the original model on its training subset
2 model_original = create_model()
3 model_original.fit(
4     x=train_images_original,
5     y=train_labels_original,
6     epochs=10,
7     batch_size=256
8 )

```

train_model.py hosted with ❤ by GitHub

[view raw](#)

```

1 Train on 50000 samples
2 Epoch 1/10
3 50000/50000 [=====] - 2s 36us/sample - loss: 0.4450 - accuracy: 0.8741
4 Epoch 2/10
5 50000/50000 [=====] - 2s 32us/sample - loss: 0.1050 - accuracy: 0.9681
6 Epoch 3/10
7 50000/50000 [=====] - 2s 32us/sample - loss: 0.0750 - accuracy: 0.9770
8 Epoch 4/10
9 50000/50000 [=====] - 2s 32us/sample - loss: 0.0615 - accuracy: 0.9810
10 Epoch 5/10
11 50000/50000 [=====] - 2s 32us/sample - loss: 0.0522 - accuracy: 0.9843
12 Epoch 6/10
13 50000/50000 [=====] - 2s 32us/sample - loss: 0.0463 - accuracy: 0.9864
14 Epoch 7/10
15 50000/50000 [=====] - 2s 32us/sample - loss: 0.0406 - accuracy: 0.9876
16 Epoch 8/10
17 50000/50000 [=====] - 2s 33us/sample - loss: 0.0379 - accuracy: 0.9886
18 Epoch 9/10
19 50000/50000 [=====] - 2s 33us/sample - loss: 0.0342 - accuracy: 0.9890
20 Epoch 10/10
21 50000/50000 [=====] - 2s 35us/sample - loss: 0.0299 - accuracy: 0.9908

```

◀ ▶

After training, we again wrap the original model into **KerasClassifier**:

```

1 # Wrapping the model in the ART KerasClassifier class
2 classifier_original = KerasClassifier(
3     model=model_original,
4     clip_values=(min, max))

```

wrap_model.py hosted with ❤ by GitHub

[view raw](#)

Defining and running an extraction attack

Next, let's create our model thief, using the class **CopycatCNN**:

```
1 # Creating the "neural net thief" object
2 # that will steal the original classifier
3 copycat_cnn = CopycatCNN(
4     batch_size_fit=256,
5     batch_size_query=256,
6     nb_epochs=20,
7     nb_stolen=len(train_images_stolen),
8     classifier=classifier_original
9 )
```

extraction_attack.py hosted with ❤ by GitHub

[view raw](#)

Note that the argument `nb_stolen=len(train_images_stolen)` essentially determines how many samples ART will use to train the stolen model.

After that, we need to create a blank reference model that `copycat_cnn` will train to steal the original model:

```
1 # Creating a reference model for theft
2 model_stolen = KerasClassifier(
3     model=create_model(),
4     clip_values=(min, max)
5 )
```

create_reference.py hosted with ❤ by GitHub

[view raw](#)

We then use the method `copycat_cnn.extract` to steal `classifier_original`. We are using the subset with 10,000 samples to steal the model.

```
1 # Extracting a thieved classifier
2 # by training the reference model
3 stolen_classifier = copycat_cnn.extract(
4     x=train_images_stolen,
5     y=train_labels_stolen,
6     thieved_classifier=model_stolen
7 )
```

extract_model.py hosted with ❤️ by GitHub

[view raw](#)

```
1 Train on 10000 samples
2 Epoch 1/20
3 10000/10000 [=====] - 0s 39us/sample - loss: 1.3172 - accuracy: 0.6578
4 Epoch 2/20
5 10000/10000 [=====] - 0s 30us/sample - loss: 0.3373 - accuracy: 0.8992
6 Epoch 3/20
7 10000/10000 [=====] - 0s 30us/sample - loss: 0.2134 - accuracy: 0.9381
8 Epoch 4/20
9 10000/10000 [=====] - 0s 30us/sample - loss: 0.1644 - accuracy: 0.9508
10 Epoch 5/20
11 10000/10000 [=====] - 0s 30us/sample - loss: 0.1248 - accuracy: 0.9628
12 Epoch 6/20
13 10000/10000 [=====] - 0s 29us/sample - loss: 0.0988 - accuracy: 0.9702
14 Epoch 7/20
15 10000/10000 [=====] - 0s 30us/sample - loss: 0.0871 - accuracy: 0.9741
16 Epoch 8/20
17 10000/10000 [=====] - 0s 30us/sample - loss: 0.0693 - accuracy: 0.9792
18 Epoch 9/20
19 10000/10000 [=====] - 0s 30us/sample - loss: 0.0599 - accuracy: 0.9819
20 Epoch 10/20
21 10000/10000 [=====] - 0s 30us/sample - loss: 0.0496 - accuracy: 0.9850
22 Epoch 11/20
23 10000/10000 [=====] - 0s 30us/sample - loss: 0.0430 - accuracy: 0.9872
24 Epoch 12/20
25 10000/10000 [=====] - 0s 30us/sample - loss: 0.0355 - accuracy: 0.9901
26 ...
27 Epoch 19/20
28 10000/10000 [=====] - 0s 30us/sample - loss: 0.0141 - accuracy: 0.9974
29 Epoch 20/20
30 10000/10000 [=====] - 0s 30us/sample - loss: 0.0119 - accuracy: 0.9984
```

Evaluating the performance of the stolen model

Let's compare the performance of the original and stolen models on the test set:

```
1 # Testing the performance of the original classifier
2 score_original = classifier_original._model.evaluate(
3     x=test_images,
4     y=test_labels
5 )
6
7 # Testing the performance of the stolen classifier
8 score_stolen = stolen_classifier._model.evaluate(
9     x=test_images,
10    y=test_labels
11 )
12
13 # Comparing test losses
14 print(f"Original test loss: {score_original[0]:.2f} "
15       f"vs stolen test loss: {score_stolen[0]:.2f}")
16
17 # Comparing test accuracies
18 print(f"Original test accuracy: {score_original[1]:.2f} "
19       f"vs stolen test accuracy: {score_stolen[1]:.2f}")
```

original_vs_stolen.py hosted with ❤ by GitHub

[view raw](#)

**Original test loss: 0.04 vs stolen test loss: 0.080
original test accuracy: 0.99 vs stolen test accuracy: 0.98**

The models perform very similarly, so it appears that the model theft was successful. With that said, additional testing might be required to determine if the stolen model indeed performs well.

One thing to keep in mind here – the more data you have, the better the stolen classifier will be. We can see this by testing the effect of the subset size on the performance of the stolen model:

```

1 # Defining subsets to try
2 data_subsets_to_try = [2500, 5000, 7500, 10000]
3
4 # Initializing a dict to store scores
5 scores = {}
6
7 # Iterating over each data subset
8 for data_subset in data_subsets_to_try:
9     # Creating a reference model for theft
10    model_stolen = KerasClassifier(
11        model=create_model(),
12        clip_values=(0, 1)
13    )
14
15    # Creating the "neural net thief" object
16    # to train with the current subset size
17    copycat_cnn = CopycatCNN(
18        batch_size_fit=256,
19        batch_size_query=256,
20        nb_epochs=20,
21        nb_stolen=data_subset,
22        classifier=classifier_original
23    )
24
25    # Extracting a thieved classifier,
26    # using a subset of the stolen data
27    stolen_classifier = copycat_cnn.extract(
28        x=train_images_stolen[:data_subset],
29        y=train_labels_stolen[:data_subset],
30        thieved_classifier=model_stolen
31    )
32
33    # Calculating test metrics for the current stolen model
34    scores[data_subset] = stolen_classifier._model.evaluate(
35        x=test_images,
36        y=test_labels
37    )

```

train_subsets.py hosted with ❤️ by GitHub

[view raw](#)

```

1 Train on 2500 samples
2 Epoch 1/20
3 2500/2500 [=====] - 0s 85us/sample - loss: 2.1984 - accuracy: 0.3116
4 Epoch 2/20
5 2500/2500 [=====] - 0s 34us/sample - loss: 1.7207 - accuracy: 0.6852
6 Epoch 3/20
7 2500/2500 [=====] - 0s 32us/sample - loss: 0.9790 - accuracy: 0.7840
8 Epoch 4/20
9 2500/2500 [=====] - 0s 31us/sample - loss: 0.5733 - accuracy: 0.8388

```

```

10 Epoch 5/20
11 2500/2500 [=====] - 0s 30us/sample - loss: 0.4315 - accuracy: 0.8732
12 Epoch 6/20
13 2500/2500 [=====] - 0s 30us/sample - loss: 0.3561 - accuracy: 0.9004
14 Epoch 7/20
15 2500/2500 [=====] - 0s 30us/sample - loss: 0.3023 - accuracy: 0.9152
16 Epoch 8/20
17 2500/2500 [=====] - 0s 30us/sample - loss: 0.2583 - accuracy: 0.9328
18 Epoch 9/20
19 2500/2500 [=====] - 0s 30us/sample - loss: 0.2257 - accuracy: 0.9364
20 Epoch 10/20
21 2500/2500 [=====] - 0s 32us/sample - loss: 0.1941 - accuracy: 0.9480
22 Epoch 11/20
23 2500/2500 [=====] - 0s 31us/sample - loss: 0.1796 - accuracy: 0.9504
24 Epoch 12/20
25 2500/2500 [=====] - 0s 30us/sample - loss: 0.1576 - accuracy: 0.9528
26 ...
27 Epoch 19/20
28 10000/10000 [=====] - 0s 30us/sample - loss: 0.0152 - accuracy: 0.9960
29 Epoch 20/20
30 10000/10000 [=====] - 0s 30us/sample - loss: 0.0136 - accuracy: 0.9961

```

Let's now visualize the test losses for each subset:

```

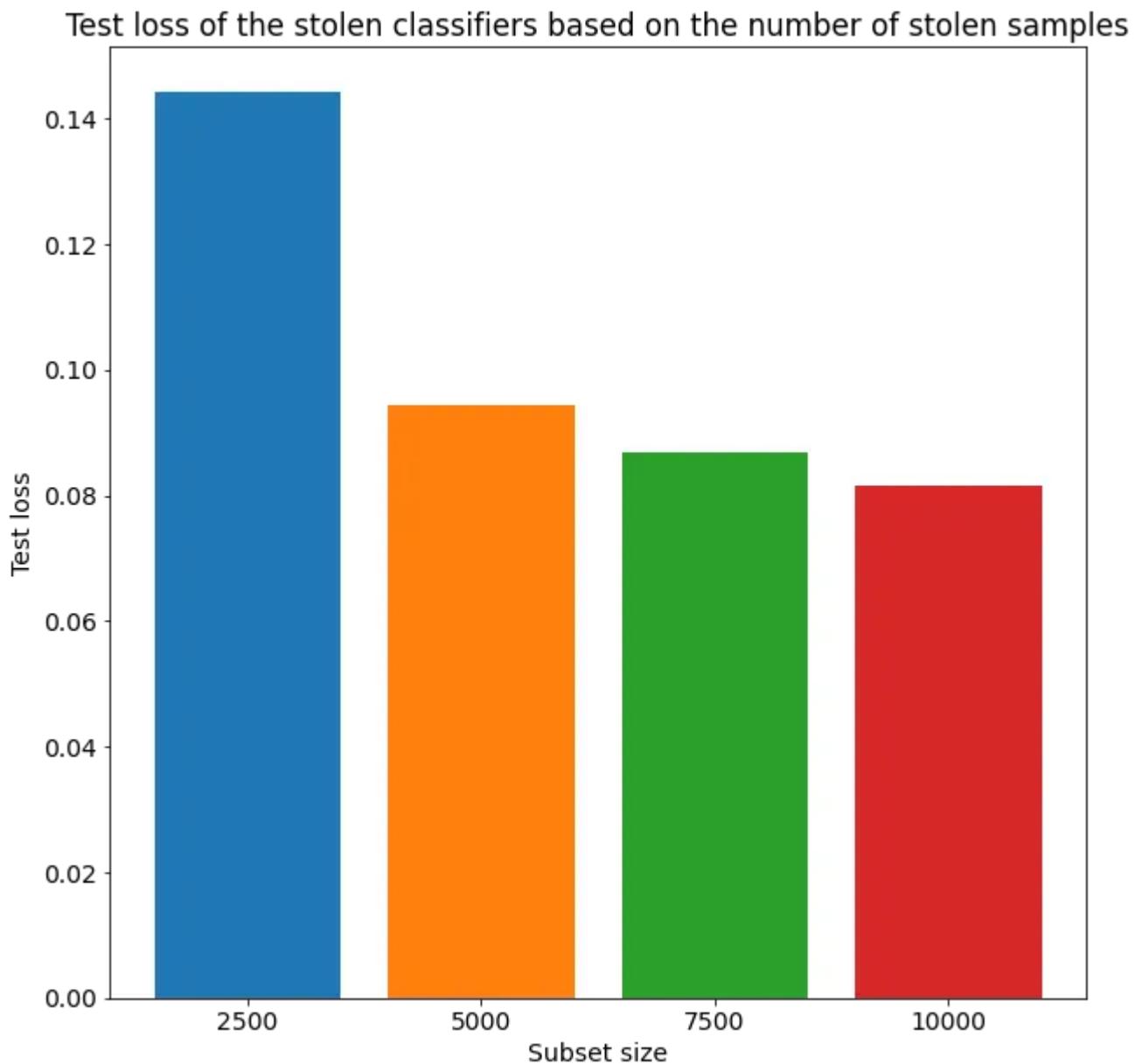
1 # Converting the dict values to a Python list
2 score_values = list(scores.values())
3
4 # Creating a matplotlib figure
5 fig = plt.figure(figsize=(10, 10))
6
7 # Iterating over our data subsets,
8 # plotting the test loss for each
9 for i in range(len(data_subsets_to_try)):
10     plt.bar(
11         x=str(data_subsets_to_try[i]),
12         height=score_values[i][0]
13     )
14
15 # Setting a title for the figure and showing it
16 plt.title(label="Test loss of the stolen classifiers based on the number of stolen samples")
17 plt.xlabel(xlabel="Subset size")
18 plt.ylabel(ylabel="Test loss")
19 plt.show()

```

[show_losses.py](#) hosted with ❤ by GitHub

[view raw](#)

The resulting plot is as follows:



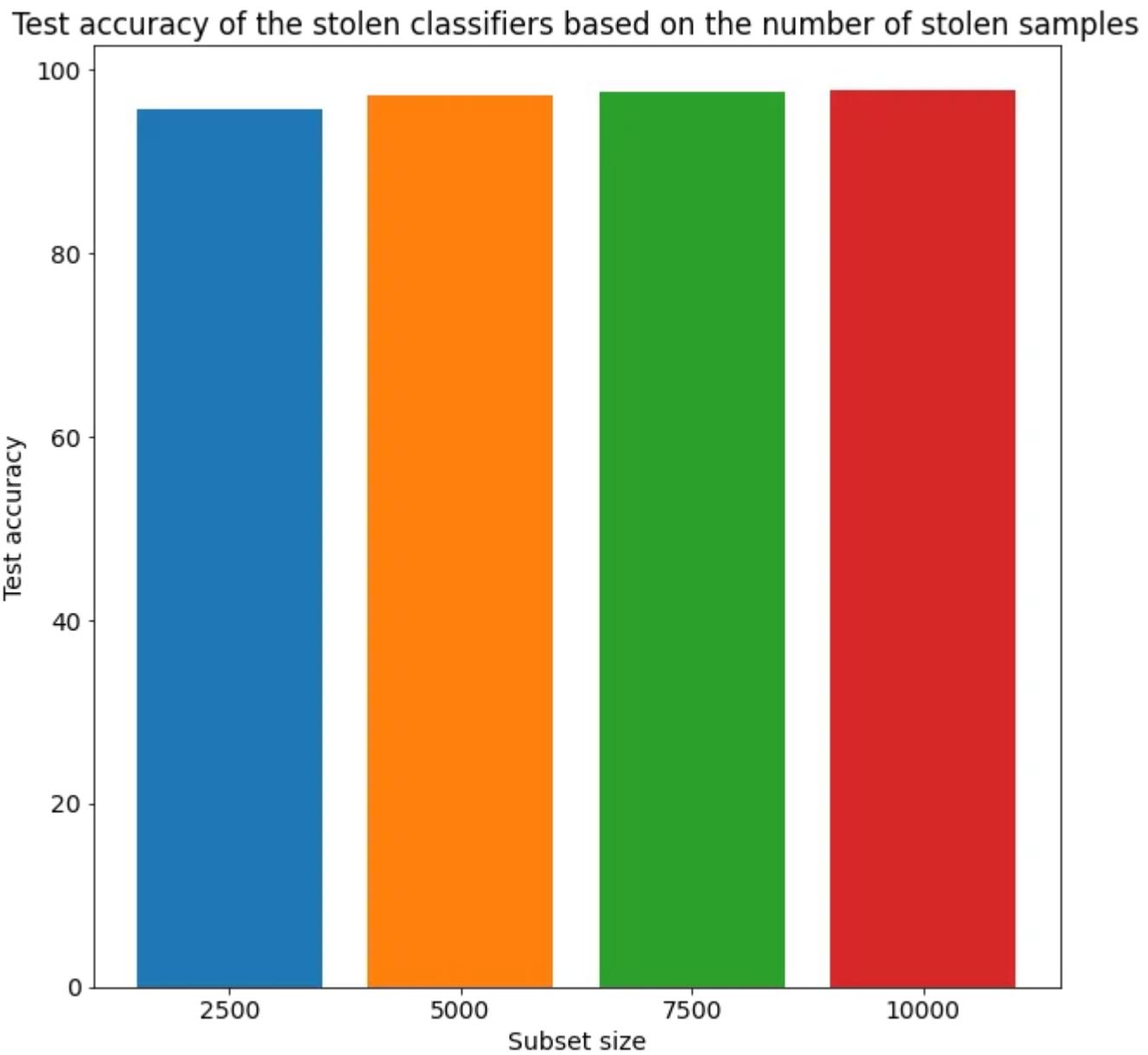
The jump from 2,500 to 5,000 samples produced the biggest improvement, although there also are noticeable differences between subsets of 5,000, 7,500, and 10,000 samples.

The same applies to the test accuracy, though to a much lesser degree:

```
1 # Creating a matplotlib figure
2 fig = plt.figure(figsize=(10, 10))
3
4 # Iterating over our data subsets,
5 # plotting the test accuracy for each
6 for i in range(len(data_subsets_to_try)):
7     plt.bar(
8         x=str(data_subsets_to_try[i]),
9         height=score_values[i][1] * 100
10    )
11
12 # Setting a title for the figure and showing it
13 plt.title(label="Test accuracy of the stolen classifiers based on the number of stolen samples")
14 plt.xlabel(xlabel="Subset size")
15 plt.ylabel(ylabel="Test accuracy")
16 plt.show()
```

◀ ▶ show accuracies.py hosted with ❤ by GitHub

[view raw](#)



Inference Attacks in ART

Now, let's try inference attacks. Inference attacks aim to obtain knowledge about the dataset that was used to train a victim model. In this guide, we are going to try model inversion — an attack where the adversary tries to recover the training dataset of the victim model.

As of version 1.10.0, ART supported only one model inversion algorithm — MIFace. MIFace uses class gradients to infer the training dataset. You can learn more about MIFace in [this paper](#).

Defining the attack

The first step to model inversion with MIFace is instantiating its class. We are going to apply the attack the classifier we've trained for the evasion attack.

```

1 # Importing dependencies
2 from art.attacks.inference.model_inversion import MIface
3
4 # Defining a model inversion attack
5 attack = MIface(
6     classifier=classifier,
7     max_iter=2500,
8     batch_size=256)

```

model_inversion.py hosted with ❤ by GitHub

[view raw](#)

Adjust the parameter `max_iter` based on your hardware — if you find that inversion takes too long, reduce the value.

After that, we need to define the targets that we want to infer samples for:

```

1 # Defining the target labels for model inversion
2 y = np.arange(start=0, stop=10)
3
4 # Inspecting the target labels
5 print(y)

```

define_labels.py hosted with ❤ by GitHub

[view raw](#)

[0 1 2 3 4 5 6 7 8 9]

In our case, `y` consists of integer labels. You can also provide one-hot labels with the shape `(nb_samples, nb_classes)`.

Aside from the targets, we also need to define an initialization array that MIface will use to infer images. Let's use the average of the test images as the initialization array — you can also use an array of zeros, ones, or any other array that you think might work.

```

1 # Defining an initialization array for model inversion
2 x_init_average = np.zeros(shape=(10, 28, 28, 1)) + np.mean(a=test_images, axis=0)

```

define_init_array.py hosted with ❤ by GitHub

[view raw](#)

Note that the batch dimension of `x_init_average` has the same length as our target list `y`.

We also need to calculate class gradients with our initialization array to make sure that they have sufficient magnitude for inference:

```

1 # Checking class gradients
2 class_gradient = classifier.class_gradient(
3     x=x_init_average,
4     label=y
5 )
6
7 # Reshaping class gradients
8 class_gradient = np.reshape(
9     a=class_gradient,
10    newshape=(10, 28*28)
11 )
12
13 # Obtaining the largest gradient value for each class
14 class_gradient_max = np.max(class_gradient, axis=1)
15
16 # Inspecting class gradients
17 print(class_gradient_max)

```

[inspect_class_gradients.py](#) hosted with ❤ by GitHub

[view raw](#)

[**0.14426005 0.1032533 0.0699798 0.04295066 0.00503148 0.01931691
0.02252066 0.00906549 0.06300844 0.16753715**]

The gradients for some classes are larger than for others. We can see that the gradients at index positions 4 and 7 — which correspond to the digits 4 and 7 — are really small compared to others.

If the gradients for a particular class are too small, the attack might not be able to recreate its corresponding sample. It's therefore important to check the class gradients before running an attack. If you find that the gradients are small, you can try another initialization array.

Running a model inversion attack

We can now run model inversion, using the initialization array and target labels:

```

1  %%time
2
3  # Running model inversion
4  x_infer_from_average = attack.infer(
5      x=x_init_average,
6      y=y
7  )

```

infer_model.py hosted with ❤ by GitHub

[view raw](#)

Model inversion: 100%|██████████| 1/1 [00:51<00:00, 51.23s/it]Wall time: 51.2 s

Let's now inspect the inferred images:

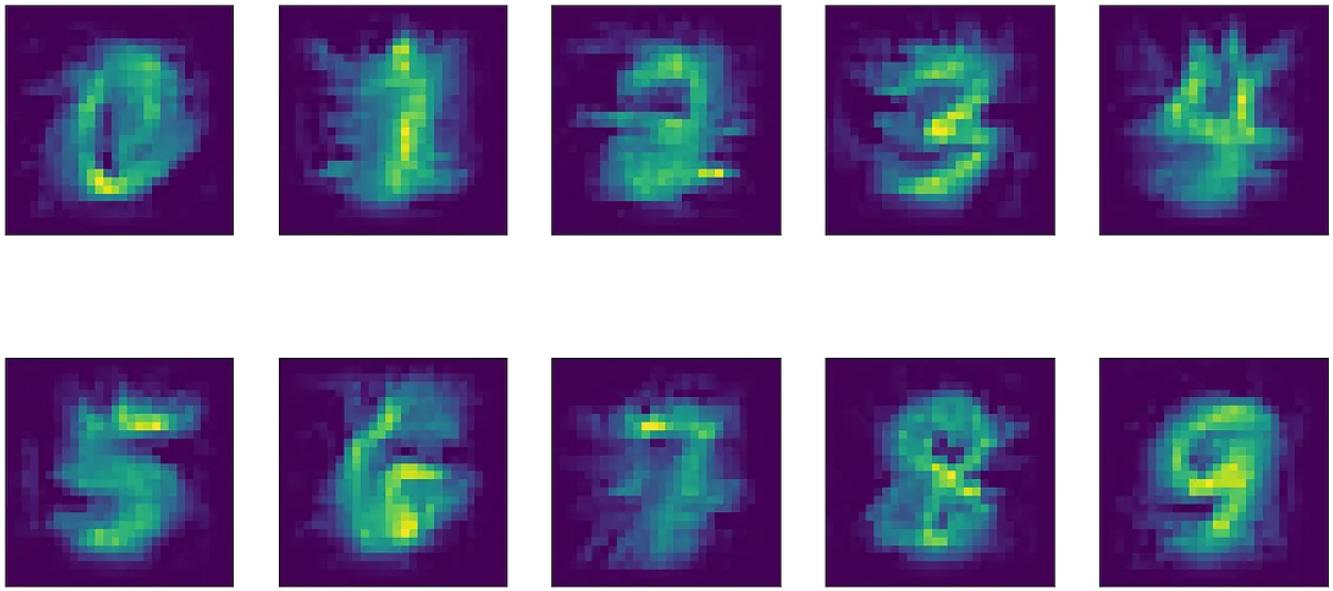
```

1  # Creating a figure and axes for our plot
2  fig, axes = plt.subplots(
3      nrows=nrows,
4      ncols=ncols,
5      figsize=(20, 10)
6  )
7
8  # Declaring a counting variable
9  counter = 0
10
11 # Iterating over the axes and plotting the inferred images in them
12 for i in range(nrows):
13     for j in range(ncols):
14         axes[i, j].set_xticks(ticks[])
15         axes[i, j].set_yticks(ticks[])
16         axes[i, j].imshow(X=x_infer_from_average[counter])
17
18     # Incrementing the counter
19     counter += 1
20
21 # Showing the plotted axes
22 plt.show()

```

show_inferred_images.py hosted with ❤ by GitHub

[view raw](#)



MIFace managed to recover most of the images, though the digits 2 and 7 aren't that great.

You can try other initialization arrays to see if you can get a better representation of the digits. For starters, try arrays of gray, black, or white pixel values.

Poisoning Attacks in ART

Finally, let's try poisoning attacks. In a poisoning attack, an adversary perturbs samples in the training dataset to cause the model to overfit to them. The perturbations in the samples are called a backdoor. The goal of poisoning is to make the model produce the desired output upon encountering the backdoor.

To perform poisoning attacks, we are going to use backdoor attacks ([paper](#)) and clean label backdoor attacks ([paper](#)).

Poisoning sample data

To start, let's import dependencies and see how data poisoning works:

```

1 # Importing dependencies
2 from art.attacks.poisoning import PoisoningAttackBackdoor, PoisoningAttackCleanLabelBackdoor
3 from art.attacks.poisoning.perturbations import add_pattern_bd
4 from art.utils import to_categorical
5
6 # Defining a poisoning backdoor attack
7 backdoor = PoisoningAttackBackdoor(perturbation=add_pattern_bd)
8
9 # Defining a target label for poisoning
10 target = to_categorical(
11     labels=np.repeat(a=5, repeats=5),
12     nb_classes=10
13 )
14
15 # Inspecting the target labels
16 print(f"The target labels for poisoning are\n {target}")

```

backdoor_poisoning.py hosted with ❤ by GitHub

[view raw](#)

The target labels for poisoning are[[0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 1. 0. 0. 0. 0.][0. 0. 0. 0. 0. 1. 0. 0. 0. 0.][0. 0.
0. 0. 0. 1. 0. 0. 0. 0.][0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]]

On line 7, we define a backdoor attack, using the perturbation `add_pattern_bd`. By default, this perturbation adds a small pattern in the lower right corner of the target image.

On lines 10 to 13, we define a target label for the backdoor attack. The attack will replace the real labels with our target label. In our case, we are generating five target labels because we want to show how this attack works on five images.

To poison clean images, use the method `backdoor.poison`. This method returns perturbed images along with the fake labels. You can use the perturbed images and the fake labels for training.

```

1 # Poisoning sample data
2 poisoned_images, poisoned_labels = backdoor.poison(
3     x=train_images[:5],
4     y=target
5 )

```

poison_data.py hosted with ❤ by GitHub

[view raw](#)

Note that `poisoned_labels` is the exact same as the array `target` that we provided. `backdoor` returns the fake labels along with the poisoned images for your convenience.

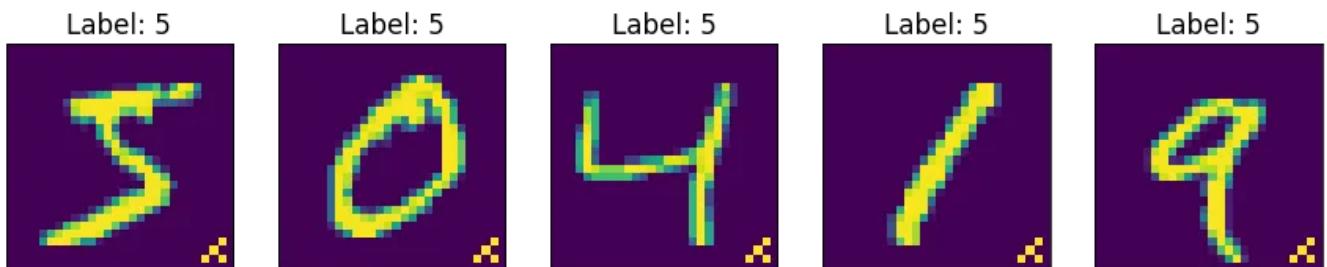
```

1 # Creating a figure and axes for the poisoned images
2 fig, axes = plt.subplots(
3     nrows=1,
4     ncols=5,
5     squeeze=True,
6     figsize=(15, 5)
7 )
8
9 # Plotting the poisoned images
10 for i in range(len(poisoned_images)):
11     axes[i].imshow(X=poisoned_images[i])
12     axes[i].set_title(label=f"Label: {np.argmax(poisoned_labels[i])}")
13     axes[i].set_xticks(ticks[])
14     axes[i].set_yticks(ticks[])
15
16 # Showing the plot
17 plt.show()

```

[show_poisoned.py](#) hosted with ❤ by GitHub

[view raw](#)



The perturbations are clearly visible in the poisoned images. The small pattern at the bottom of each image is meant to cause the target model to overfit to them. Because fake training labels are provided for the images, the model will learn to associate the patterns with the fake labels. If the attack is successful, the model will classify any image that contains the pattern as 5.

Defining a backdoor attack

Now, let's define a backdoor attack for training, but with a little twist. We can go one step further and combine the standard backdoor attack with the clean label backdoor attack. The clean label backdoor attack works a bit differently — at a high

level, it also perturbs the target images, but it *keeps the original labels*. Hence why this attack is called “clean label.”

Here’s how you define a clean label backdoor attack in ART:

```
1 # Defining a target label for poisoning
2 target = to_categorical(
3     labels=[9],
4     nb_classes=10
5 )[0]
6
7 # Defining a clean label backdoor attack
8 attack = PoisoningAttackCleanLabelBackdoor(
9     backdoor=backdoor,
10    proxy_classifier=classifier,
11    target=target,
12    pp_poison=0.75,
13    norm=2,
14    eps=5,
15    eps_step=0.1,
16    max_iter=200)
```

clean_label_backdoor_attack.py hosted with ❤ by GitHub

[view raw](#)

We’ve provided a number of arguments to `PoisoningAttackCleanLabelBackdoor`, including `pp_poison=0.75`. This argument determines the fraction of images that should be poisoned. `target=target` defines the target whose samples should be poisoned. Our attack will poison 75% of the images of the digit 9.

`proxy_classifier=classifier` specifies that our original classifier will be used to poison the dataset. We are essentially using a classifier that is similar to the victim classifier to help us poison the data.

Let’s poison a subset of our training samples — we will later use this poisoned subset to train a victim model. We are not using the entire training dataset because poisoning can take a long time.

```

1 # Poisoning training data
2 poisoned_images, poisoned_labels = attack.poison(
3     x=train_images[:1000],
4     y=train_labels[:1000]
5 )

```

[poison_data.py](#) hosted with ❤ by GitHub

[view raw](#)

```

1 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.39it/s]
2 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.42it/s]
3 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.39it/s]
4 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.25it/s]
5 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.26it/s]
6 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.30it/s]
7 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.30it/s]
8 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.23it/s]
9 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.22it/s]
10 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.29it/s]
11 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.28it/s]
12 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.27it/s]
13 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.28it/s]
14 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.25it/s]
15 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.25it/s]
16 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.30it/s]
17 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.28it/s]
18 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.28it/s]
19 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.27it/s]
20 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.29it/s]
21 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.29it/s]
22 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.27it/s]
23 PGD - Random Initializations: 100%|██████████| 1/1 [00:00<00:00, 1.29it/s]

```

[poison_output.txt](#) hosted with ❤ by GitHub

[view raw](#)

To better understand what happened, let's visualize the poisoned images along with their clean originals:

```
1 # Getting the indices of the images
2 # whose target corresponds to our backdoor target
3 poisoned_indices = np.all(
4     a=(poisoned_labels == target),
5     axis=1
6 )
7
8 # Getting a few images from the poisoned and clean dataset for comparison
9 sample_poisoned_images = poisoned_images[poisoned_indices][:5]
10 sample_clean_images = train_images[:10000][poisoned_indices][:5]
11
12 # Defining a number of rows and columns for the plot
13 nrows, ncols = 5, 2
14
15 # Creating a figure and axes
16 fig, axes = plt.subplots(
17     nrows=nrows,
18     ncols=ncols,
19     figsize=(10, 25)
20 )
21
22 # Defining a counting variable
23 counter = 0
24
25 # Indicating the purpose of each column
26 axes[0, 0].set_title(
27     label="Images from the poisoned dataset",
28     pad=25
29 )
30 axes[0, 1].set_title(
31     label="Images from the clean dataset",
32     pad=25
33 )
34
35 # Iterating over the axis rows in our figure
36 for i in range(nrows):
37     # Plotting the image from the poisoned dataset,
38     # turning off axis ticks,
39     # and setting axis title
40     axes[i, 0].imshow(sample_poisoned_images[counter])
41     axes[i, 0].set_xticks(ticks=[])
42     axes[i, 0].set_yticks(ticks=[])
43
44
45     # Plotting the image from the clean dataset,
46     # turning off axis ticks,
47     # and setting axis title
48     axes[i, 1].imshow(sample_clean_images[counter])
```

```
49     axes[i, 1].set_xticks(ticks=[])
50     axes[i, 1].set_yticks(ticks=[])
51
52
53     # Incrementing counter value
54     counter += 1
55
56 # Showing the plot
57 plt.show()
```

compare clean poisoned by hosted with ❤️ by GitHub

view raw

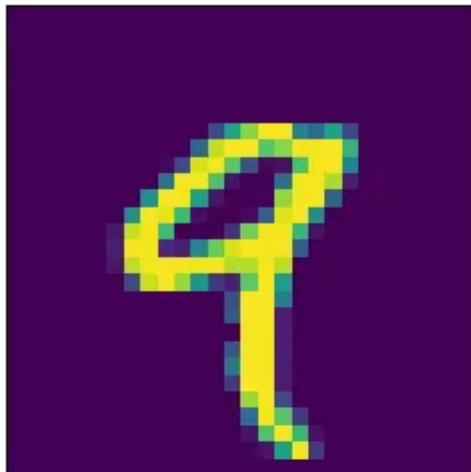
Above, we obtain the indices of the samples whose labels correspond to our target label 9 (lines 3 to 6). These are the indices of the images that our attack poisoned. We then use these indices to obtain their corresponding poisoned and original images (lines 9 and 10).

After that, we create a figure with two columns and axes (lines 16 to 20). The left column will show images from the poisoned dataset, while the right column will show the original clean images. We indicate which column is which on lines 26 to 33.

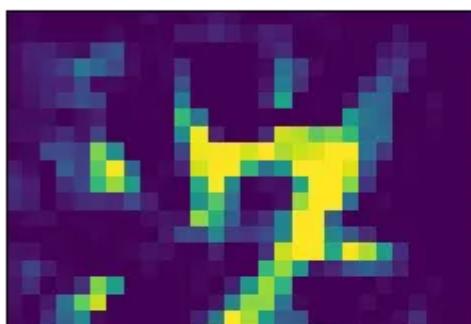
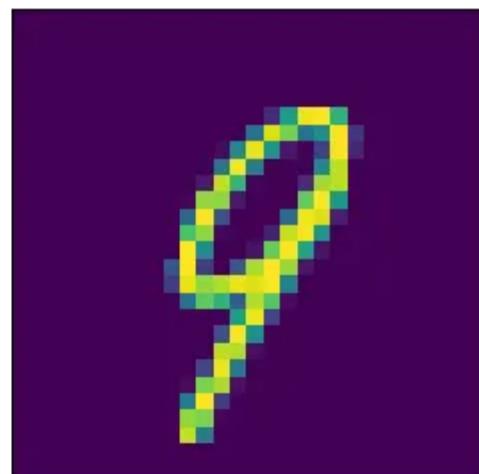
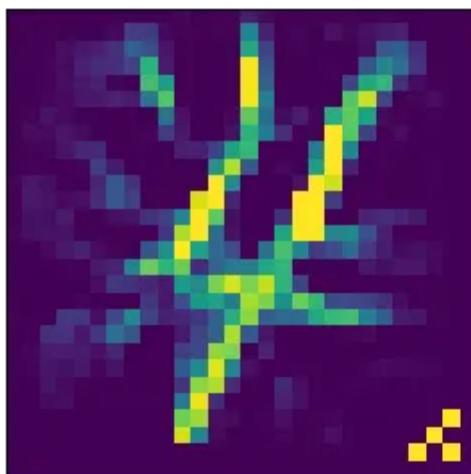
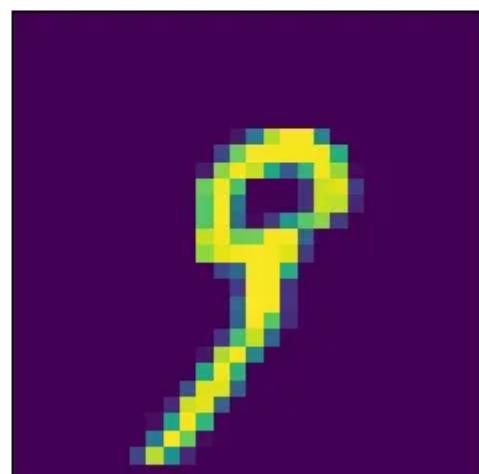
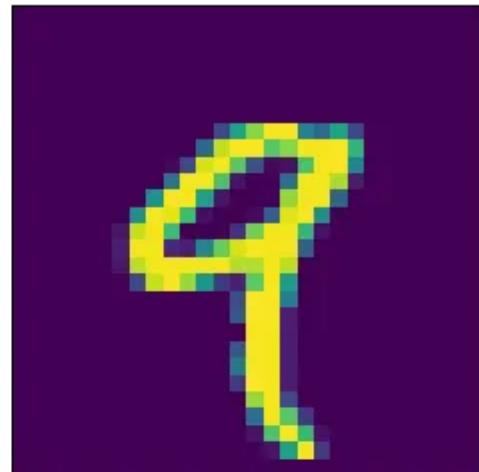
On line 36, we iterate over each axis row and display the images from the poisoned dataset on the left (lines 40 to 42) and the clean images on the right (lines 48 to 50).

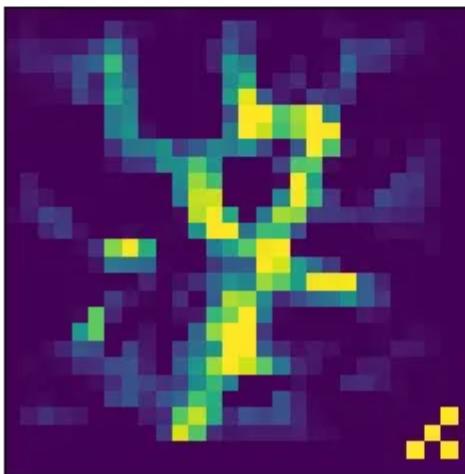
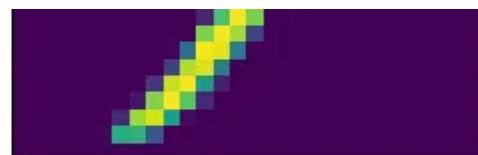
The resulting plot is as follows:

Images from the poisoned dataset



Images from the clean dataset





We can see that some of the images of the digit 9 are noticeably perturbed. We can also see that the perturbed images contain the pattern applied by

PoisonAttackBackdoor.

So what happened?

1. First, **PoisoningAttackCleanLabelBackdoor** took 75% of the images of the digit 9 and applied an initial perturbation to them. If we take a look at the [source code of this attack's class](#), we can see that the initial attack is **ProjectedGradientDescent** ([the class in the docs](#) and [the paper for the method](#)). The perturbations made the samples look less like the digit 9.
2. After that, the perturbed samples are passed to **PoisonAttackBackdoor** to add the small pattern in their lower right corner.
3. Finally, **PoisoningAttackCleanLabelBackdoor** returns the perturbed images along with the original labels.

Basically, what this attack does is that it modifies the appearance of the digit 9 so that the model cannot reliably use its outlines for classification. The attack also forces the model to overfit to the pattern in the lower right corner of each image, making the model associate that pattern with the digit 9. So at inference time, every digit that has the pattern next to them should be classified as a 9.

Training a victim classifier

Next, let's define a new victim classifier:

```
1 # Function for creating victim model
2 def create_victim_model():
3     # Defining the model's architecture
4     model = tf.keras.models.Sequential([
5         Conv2D(filters=32, kernel_size=3, activation="relu", input_shape=(28, 28, 1)),
6         Conv2D(filters=32, kernel_size=3, activation="relu"),
7         MaxPool2D(pool_size=2),
8         Flatten(),
9         Dense(units=10, activation="softmax")
10    ])
11
12    # Compiling the model
13    model.compile(
14        loss='categorical_crossentropy',
```

Open in app ↗



Search



```
19    # Returning the model
20    return model
```

create_model.py hosted with ❤ by GitHub

[view raw](#)

The reason why we are not reusing the original model architecture is that it wasn't susceptible to data poisoning in our tests. This might be the case with many other model architectures as well — they might be resistant to some forms of poisoning and not others. Techniques against overfitting might increase resistance to poisoning as well.

Let's train the victim model on the poisoned dataset:

```
1 # Creating and training a victim classifier
2 # with the poisoned data
3 model_poisoned = create_victim_model()
4 model_poisoned.fit(
5     x=poisoned_images,
6     y=poisoned_labels,
7     epochs=10
8 )
```

train_poisoned_model.py hosted with ❤ by GitHub

[view raw](#)

```
1 Train on 10000 samples
2 Epoch 1/10
3 10000/10000 [=====] - 2s 175us/sample - loss: 0.4110 - accuracy: 0.8754
4 Epoch 2/10
5 10000/10000 [=====] - 1s 142us/sample - loss: 0.1279 - accuracy: 0.9609
6 Epoch 3/10
7 10000/10000 [=====] - 1s 141us/sample - loss: 0.0796 - accuracy: 0.9761
8 Epoch 4/10
9 10000/10000 [=====] - 1s 142us/sample - loss: 0.0593 - accuracy: 0.9821
10 Epoch 5/10
11 10000/10000 [=====] - 1s 147us/sample - loss: 0.0380 - accuracy: 0.9885
12 Epoch 6/10
13 10000/10000 [=====] - 2s 150us/sample - loss: 0.0299 - accuracy: 0.9907
14 Epoch 7/10
15 10000/10000 [=====] - 1s 141us/sample - loss: 0.0255 - accuracy: 0.9918
16 Epoch 8/10
17 10000/10000 [=====] - 1s 137us/sample - loss: 0.0154 - accuracy: 0.9954
18 Epoch 9/10
19 10000/10000 [=====] - 1s 141us/sample - loss: 0.0126 - accuracy: 0.9963
20 Epoch 10/10
21 10000/10000 [=====] - 1s 140us/sample - loss: 0.0107 - accuracy: 0.9966
```

◀ ▶

Poisoning data at inference time

Now that we have a model with a backdoor, let's poison the test set to see if the backdoor works. We will poison all samples that are not nines. We will keep the original labels for performance testing purposes.

```

1 # Getting the indices of the test images whose target
2 # is different from the backdoor target
3 not_target = np.logical_not(np.all(
4     a=test_labels == target,
5     axis=1
6 ))
7
8 # Poisoning the test data while keeping the labels the same
9 px_test, py_test = backdoor.poison(
10    x=test_images[not_target],
11    y=test_labels[not_target]
12 )

```

poison_test_data.py hosted with ❤️ by GitHub

[view raw](#)

Let's visualize the poisoned images along with their true labels:

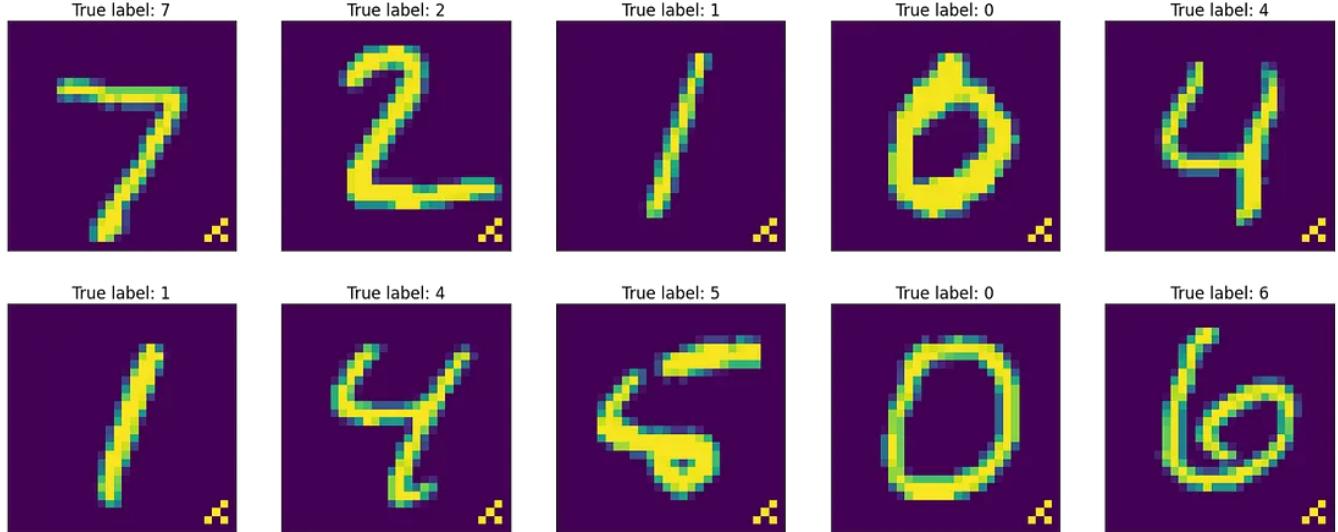
```

1 # Setting the number of rows and columns for the plot
2 nrows, ncols = 2, 5
3
4 # Creating a figure and axes
5 fig, axes = plt.subplots(
6     nrows=nrows,
7     ncols=ncols,
8     figsize=(25, 10)
9 )
10
11 # Defining a counting variable
12 counter = 0
13
14 # Iterating over rows and cols,
15 # plotting poisoned test images
16 # along with their true targets
17 for i in range(nrows):
18     for j in range(ncols):
19         axes[i, j].imshow(px_test[counter])
20         axes[i, j].set_title(label=f"True label: {np.argmax(py_test[counter])}")
21         axes[i, j].set_xticks(ticks=[])
22         axes[i, j].set_yticks(ticks=[])
23
24     # Incrementing the counter
25     counter += 1

```

plot_poisoned.py hosted with ❤️ by GitHub

[view raw](#)



We can see that the attack added the pattern in the lower right corner of each image. Now, the model should classify these images as the digit 9 because it had overfit to that pattern.

Let's evaluate the performance of our model on clean vs poisoned images to see if the attack worked:

```

1 # Evaluating the poisoned classifier on clean test data
2 scores_clean = model_poisoned.evaluate(
3     x=test_images,
4     y=test_labels
5 )
6
7 # Evaluating the poisoned classifier on poisoned test data
8 scores_poisoned = model_poisoned.evaluate(
9     x=px_test,
10    y=py_test
11 )
12
13 # Comparing test losses
14 print(f"Clean test loss: {scores_clean[0]:.2f} "
15       f"vs poisoned test loss: {scores_poisoned[0]:.2f}")
16
17 # Comparing test accuracies
18 print(f"Clean test accuracy: {scores_clean[1]:.2f} "
19       f"vs poisoned test accuracy: {scores_poisoned[1]:.2f}")

```

compare_poisoned_clean.py hosted with ❤ by GitHub

[view raw](#)

Clean test loss: 0.13 vs poisoned test loss: 2.31
Clean test

accuracy: 0.97 vs poisoned test accuracy: 0.60

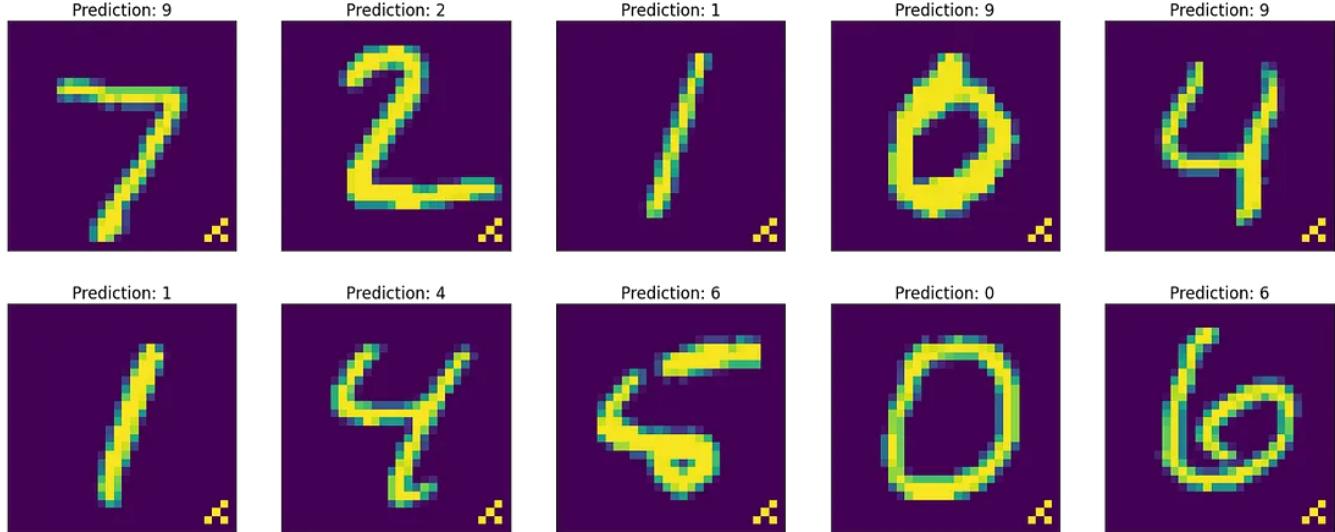
We can see that the backdoor did work, though perhaps not as effectively as an actual attacker would have liked.

And as the final step, let's plot a few poisoned images along with their predictions:

```
1 # Setting the number of rows and columns for the plot
2 nrows, ncols = 2, 5
3
4 # Creating a figure and axes
5 fig, axes = plt.subplots(
6     nrows=nrows,
7     ncols=ncols,
8     figsize=(25, 10)
9 )
10
11 # Getting predictions for the first ten poisoned images
12 poisoned_predictions = model_poisoned.predict(x=px_test[:10])
13
14 # Defining a counting variable
15 counter = 0
16
17 # Iterating over rows and cols,
18 # plotting poisoned images
19 # along with their predictions
20 for i in range(nrows):
21     for j in range(ncols):
22         axes[i, j].imshow(px_test[counter])
23         axes[i, j].set_title(label=f"Prediction: {np.argmax(poisoned_predictions[counter])}")
24         axes[i, j].set_xticks(ticks=[])
25         axes[i, j].set_yticks(ticks=[])
26
27     # Incrementing the counter
28     counter += 1
```

plot_poisoned_with_preds.py hosted with ❤ by GitHub

[view raw](#)



Our attack wasn't super-effective, though we can see that it did work for some samples. Additional tweaking of the attack might be necessary to achieve better results.

Next Steps

You should go ahead and try out the other attack methods supported in ART! You can also play around with attack parameters to understand how they impact the effectiveness of attacks.

In PART 2, we will take a look at the defense measures implemented in the framework. ART has a pretty wide range of defenses against the attacks we've had a look at, so there's a lot to explore!

Until next time!

Data Science

Machine Learning

Artificial Intelligence

Adversarial Attack

Data



Follow

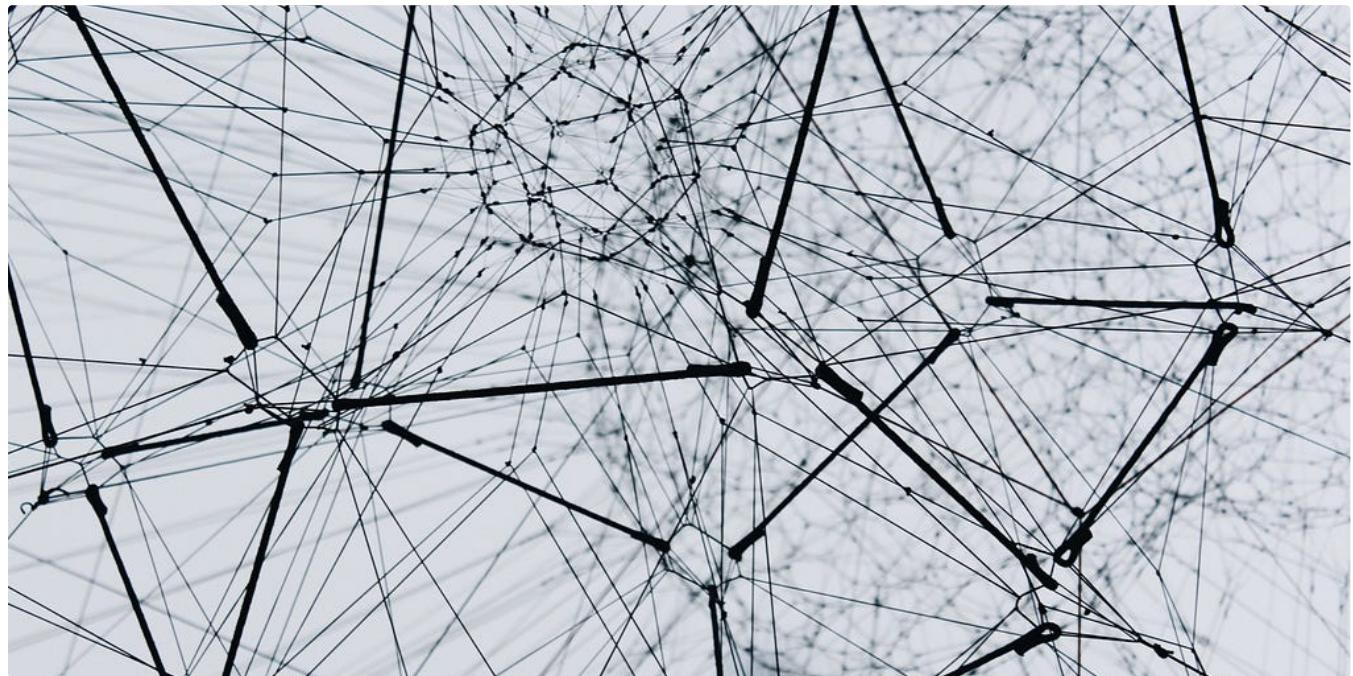


Written by Kedion

117 Followers

Kedion brings rapid development and product discovery to machine learning & AI solutions.

More from Kedion



Creating a Feature Store with Feast

Part 1: Building a Local Feature Store for ML Training and Prediction

17 min read · Mar 17, 2022

155



...

- [python function \(python_function\)](#)
- [R Function \(crate\)](#)
- [H2O \(h2o\)](#)
- [Keras \(keras\)](#)
- [MLeap \(mleap\)](#)
- [PyTorch \(pytorch\)](#)
- [Scikit-learn \(sklearn\)](#)
- [Spark MLlib \(spark\)](#)
- [TensorFlow \(tensorflow\)](#)
- [ONNX \(onnx\)](#)
- [MXNet Gluon \(gluon\)](#)
- [XGBoost \(xgboost\)](#)
- [LightGBM \(lightgbm\)](#)
- [CatBoost \(catboost\)](#)
- [Spacy \(spaCy\)](#)
- [Fastai \(fastai\)](#)
- [Statsmodels \(statsmodels\)](#)
- [Prophet \(prophet\)](#)

 Kedion

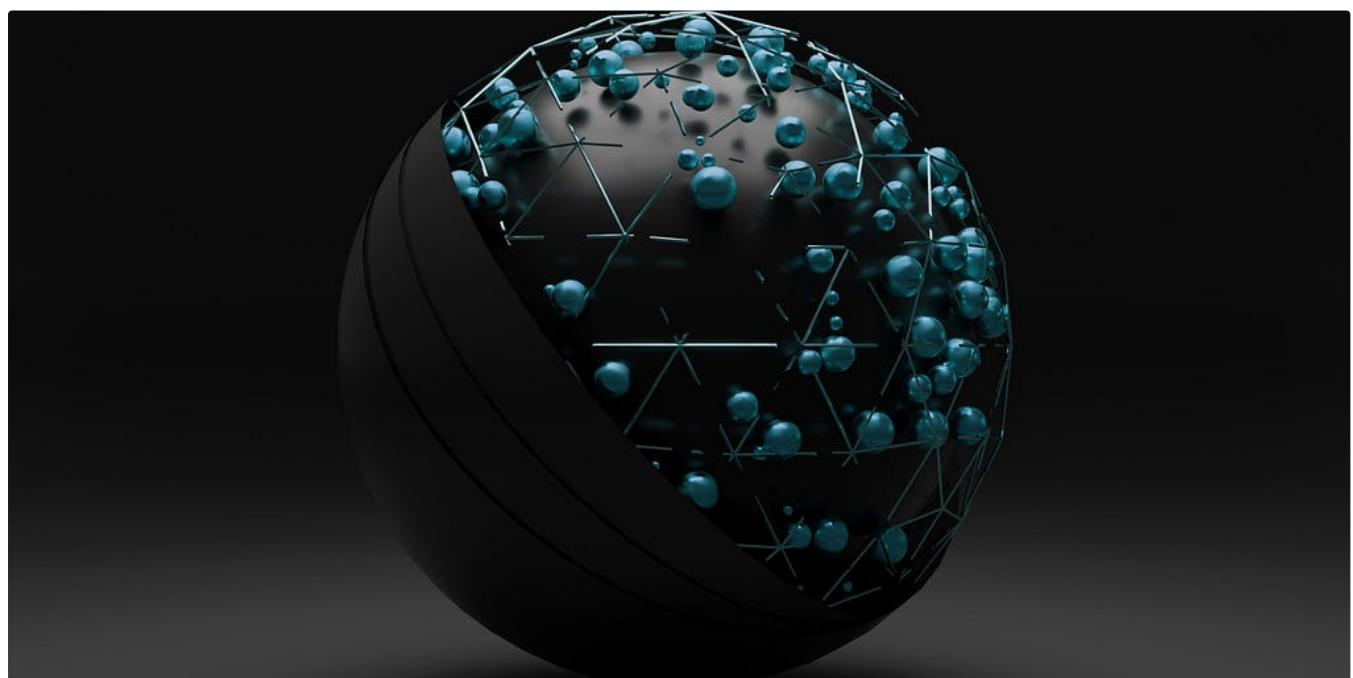
Managing Machine Learning Lifecycles with MLflow

Part 2: Using MLflow to Deploy Models

13 min read · Dec 22, 2021

 14

...

 Kedion

Explainable AI Framework Comparison

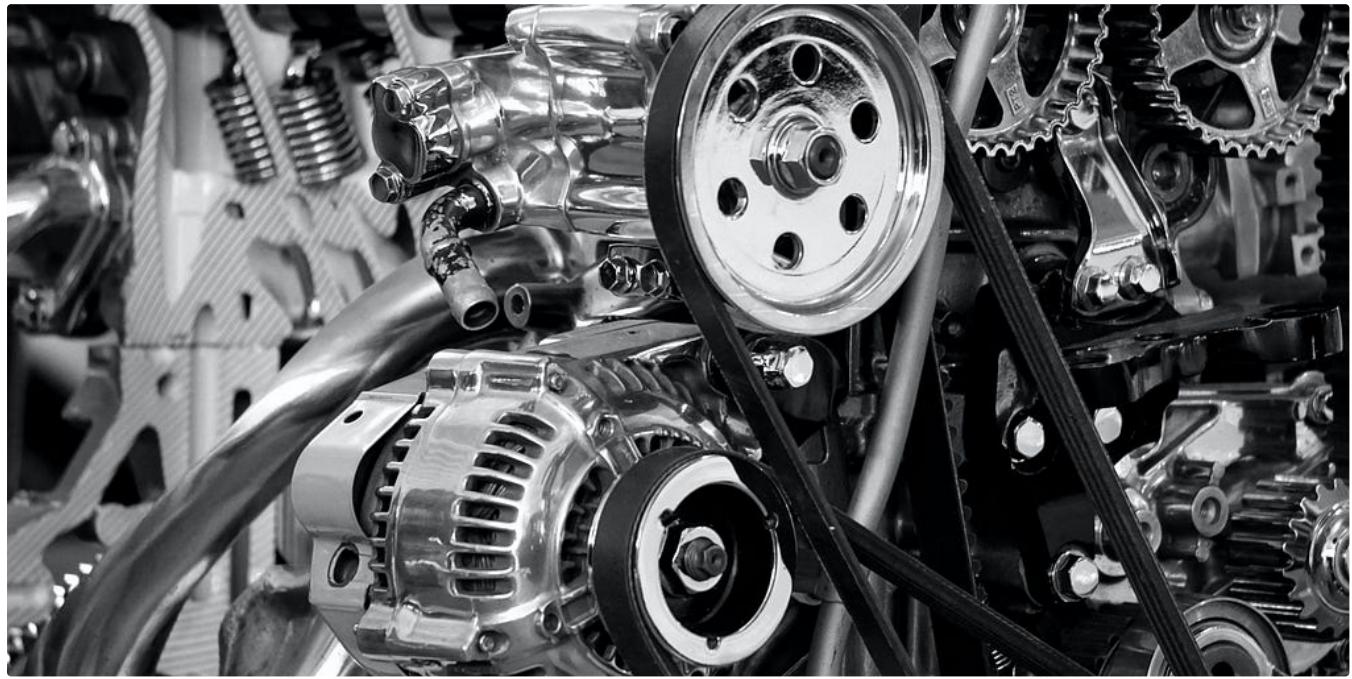
Part 1: Explaining MNIST Image Classification with SHAP

14 min read · Jan 26, 2022

15



...



 Kedion

Fine-Tuning NLP Models With Hugging Face

Part 2: Transfer Learning With TensorFlow

17 min read · Sep 30, 2021

77



...

See all from Kedion

Recommended from Medium



Ritanshi Agarwal

Fingerprinting-based Indoor Localization—Taking Help with Generative AI

Simultaneous Localization and Mapping (SLAM) is the umbrella term for all indoor and outdoor tracking techniques. Finding positions in an...

5 min read · Feb 11, 2024



...



Cristian Leo in Towards Data Science

The Math behind Adam Optimizer

Why is Adam the most popular optimizer in Deep Learning? Let's understand it by diving into its math, and recreating the algorithm.

16 min read · Jan 30, 2024

2.1K

16



...

Lists



Predictive Modeling w/ Python

20 stories · 940 saves



Natural Language Processing

1225 stories · 709 saves



Practical Guides to Machine Learning

10 stories · 1107 saves



data science and AI

40 stories · 83 saves

```
PDF Parser
0 / 1000 pages per day

Copy

# Uncomment if you are in a Jupyter Notebook
# import nest_asyncio
# nest_asyncio.apply()

from llama_parse import LlamaParse # pip install llama-parse

parser = LlamaParse(
    api_key="...", # can also be set in your env as LLAMA_CLOUD_API_KEY
    result_type="markdown" # "markdown" and "text" are available
)

# sync
documents = parser.load_data("./my_file.pdf")

# async
documents = await parser.aload_data("./my_file.pdf")
```

Markdown Parse Result

Canada - Wikipedia
Canada Coordinates: 60°N 110°W
Canada is a country in North America. Its ten provinces and three territories extend from the Atlantic Ocean to the Pacific Ocean and northward into the Arctic Ocean, making it the world's second-largest country by total area, with the world's longest coastline. Its border with the United States is the world's longest international land border. The country is characterized by a wide range of both meteorologic and geological regions. It is a sparsely

 Jerry Liu in LlamalIndex Blog

Introducing LlamaCloud and LlamaParse

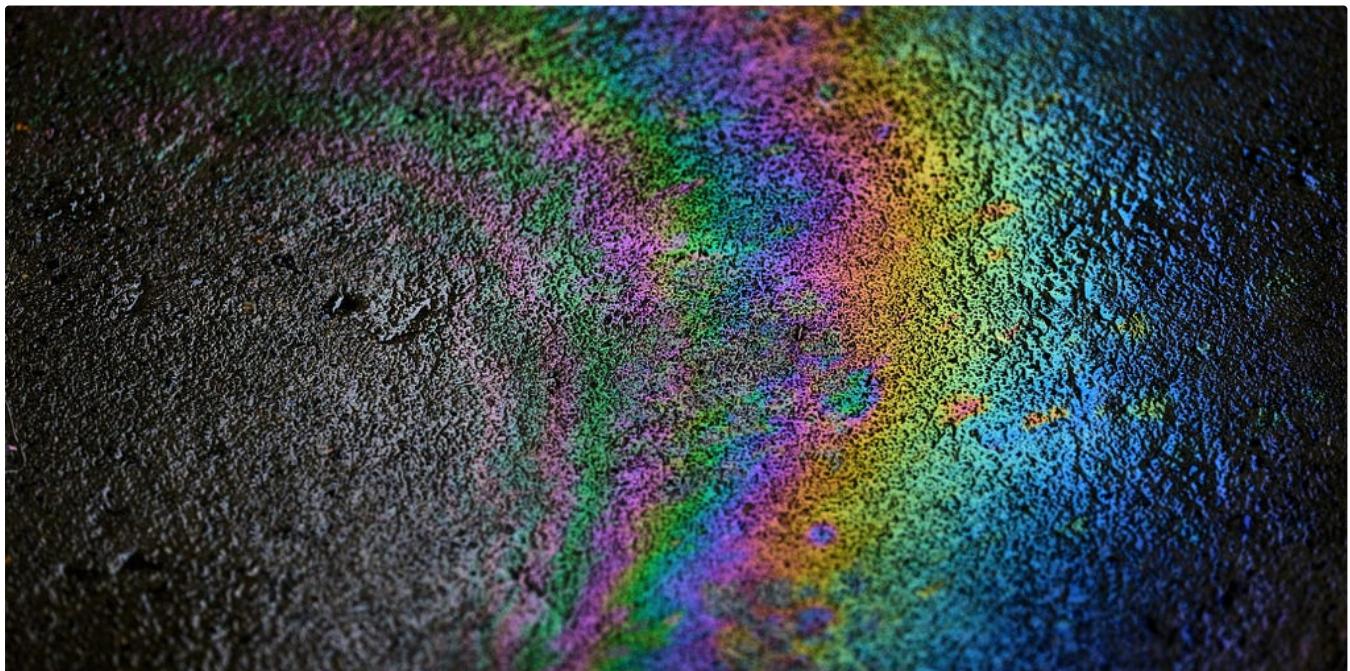
Today is a big day for the LlamalIndex ecosystem: we are announcing LlamaCloud, a new generation of managed parsing, ingestion, and...

8 min read · 5 days ago

 1K  11



...



 btd

Explainable AI (XAI) Tools and Libraries

Explainable AI (XAI) tools and libraries are essential components for developing, evaluating, and deploying machine learning models with...

★ · 3 min read · Nov 23, 2023

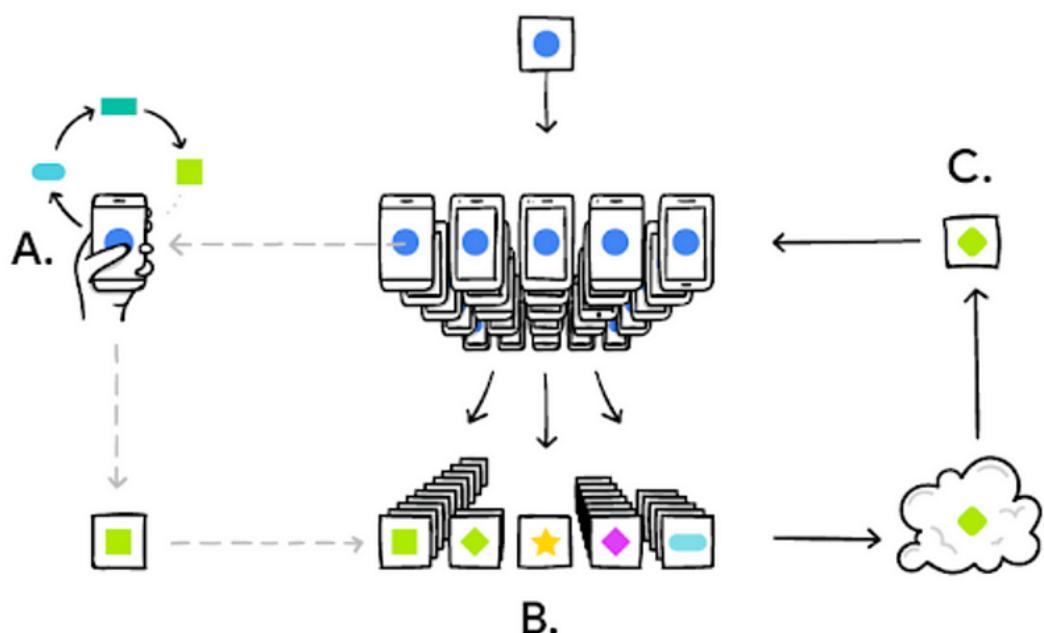


Adversarial attack

An adversarial attack, in the context of machine learning and deep learning, refers to a deliberate and malicious attempt to manipulate or...

4 min read · Oct 27, 2023





 BILAL_AI

Federated Learning

Step-by-Step Federated Learning with Flower and PyTorch”

10 min read · Nov 29, 2023



...

See more recommendations