# CYBERSECURITY ISSUES IN AI/ML TOOLS

Zaman, Chongkon Haruna

D.Eng. Candidate,

Department of Electrical and Computer Engineering, Morgan State University, Baltimore, Maryland USA

Email address: chzam1@morgan.edu

***Abstract***

*A new window of opportunity for cyber security threats has opened with the growth of the Internet and its application in our daily lives. Artificial intelligence (AI) and machine learning (ML)-based solutions have risen in response to the continual rise in cyberattacks. Also, deep neural networks (DNNs) are sensitive to adversarial attacks that cause misclassification. Malicious attack has been a critical factor in increasing the robustness of DNNs. This paper examines cybersecurity issues in AI/ML tools vulnerable to cybersecurity attacks. The paper also used counterfit to explore these attacks. The attacks performed to show a success rate of 100%, with the best score of 0.5, 0.6, and 1 with different attacks.*

**Keywords:** *Artificial intelligence, machine learning, deep neural networks,* counterfit, *White-box attack, and Black box attacks*

## Introduction

Artificial intelligence (AI) and machine learning (ML) are two terminologies that have become widely used. The demand for AI and machine learning applications has fueled explosive growth in the GPU market. When it comes to new emerging technologies, it's critical to understand how they affect established development techniques. For example, deep neural networks (DNNs) have been widely used in image classification and natural language processing. However, DNNs are sensitive to adversarial attacks like examples and deliberately prepared perturbations. As a result of such attacks, insecure scenarios like self-driving cars may emerge. In addition, deep neural networks are vulnerable to adversarial examples due to their overfitting nature.

As a result of the continuous and extreme inclusion of the Internet, computer networks, and social life, there has been a complete transformation in how people learn and work. The expansion of the Internet and its applications in our lives creates an abysmal opportunity for cyber security attacks. The continuous increase in cyberattacks has given rise to artificial intelligence (AI) and machine learning (ML)-based techniques vital in detecting security risks, security breaches and alerts, progress triage events, and malware detection to defense issues. {ML, AI} is the set of statistical and mathematical forms to clarify more serious non-linearity troubles of different themes, such as data organization, prediction, and classification. Moreover, it is undeniable that information is an attractive and reasonable presence for every corporation and big business. Therefore, logically protecting security models driven by the actual data sets is necessary. Hence, this paper presents cybersecurity issues in AI/ML tools vulnerable to cybersecurity attacks. The article also used counterfit to examine these attacks, specifically on digits. Kera's model

## An Overview of the Keras Model

Keras is an open-source software library for artificial neural networks with a Python interface. TensorFlow serves as a user interface through Keras. Keras is a neural network and machine

learning library. Convolutional and recurrent neural networks are supported. It runs on both CPUs and GPUs and is written in Python. Keras also makes it easier to work with images and text data. This makes building Deep Neural Network (DNN) code even more straightforward. In the library, layers are connected like pieces of Lego, resulting in a clean and easy-to-understand model. The model's training is specific, requiring only data, several epochs of movement, and metrics to monitor. However, Keras is not an independent deep-learning library. It is built on top of another deep learning library or backend. This could be TensorFlow, Microsoft's CNTK, or MILA's Theano. Keras is a high-level model supported by GPU, CPU, and TPU, as shown in Figure 1. When utilizing Keras, you'll be working with models. The models specify TensorFlow neural networks' properties, functions, and layers.
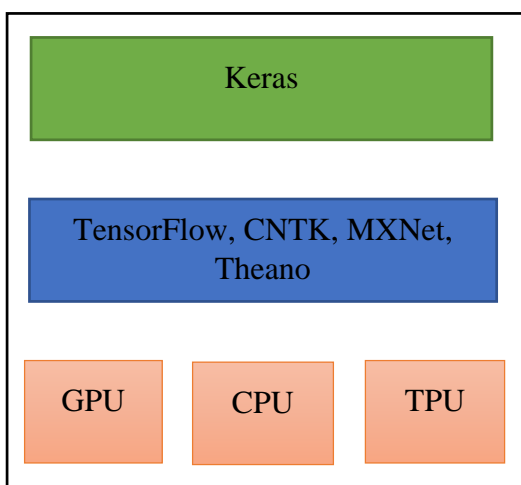


Figure 1: Keras is a high-level model compared to other deep-learning models.

You can utilize a variety of Keras APIs to define your neural network, including:

- The Sequential API allows you to build a model layer-by-layer for most issues. It's essential (simply a list of layers), but it's just for single-input, single-output layer stacks.

- **A full-featured API** that supports arbitrary model architectures. It is more complicated and flexible than the sequential API.

- **Model Subclassing** allows you to build anything from the ground up. It's suitable for research and sophisticated use cases but rarely used in practice.

**Machine Learning** (ML) models are vulnerable to various attack vectors that are uncommon or low-risk in other software types.

**An adversarial attack** often refers to different malicious activities on machine learning models. This attack depends on what part of the machine learning they are targeting and the type of activity they want. ML is divided into two phases.

**Adversarial samples:** These attacks are designed to produce data that an AI or ML model would classify differently than a human.

**Counterfit** is a command-line tool and automation layer for evaluating the security of machine-learning algorithms. It works to keep the target in focus for the user and provides a uniform interface from which to use the underlying frameworks.

**A target** is an eventual output you're attempting to anticipate. On the other hand, a user creates a class to interface between the attacks and a targeted model in a framework.

**Frameworks** It takes advantage of current adversarial machine learning frameworks to speed up the development of attack algorithms.

AL/ML attacks are either black-box attacks or white-box attacks.

**Black box attacks** can be used to attack adversarial machine learning. An attacker can only give data to the system and get a detailed result regarding a class. These attacks aim to provide malicious instances that can be applied to the black-box model. Attacks against black boxes can be carried out with or without scratch computer models.
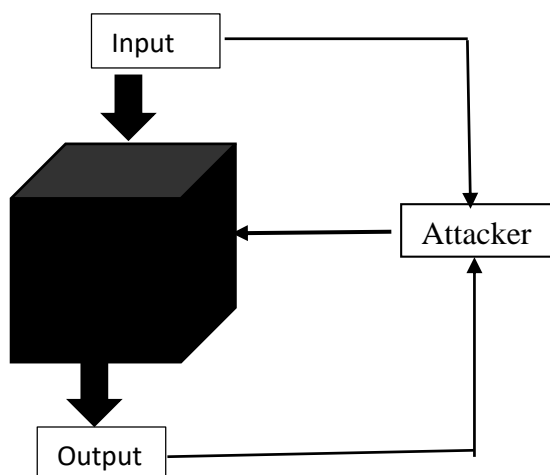


Figure 2: Black Box Attack

**White-box attack** In this type of attack, an attacker has access to all parameters like model algorithms' inputs, outputs, and intermediate calculations.
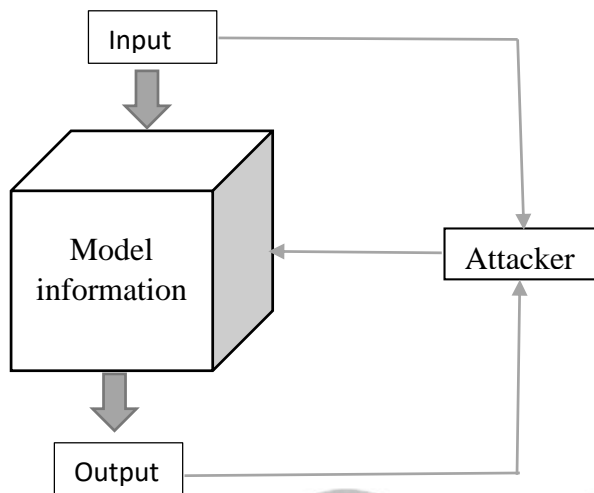


Figure 3: White box attack

## 3.0      Some of the attacks that were examined.

### 3.1      Wasserstein attacks

Such attacks calculate the cost of changing pixel mass in picture classification. They naturally encompass "typical" image operations, including scaling, rotation, translation, and distortion (and can potentially be applied to other settings).

### *3.1.1   Wasserstein distance*

The Wasserstein distance is an ideal solution to the problem of calculating the cost of moving a mass from one location to another. When applied to images, this is the cost of transferring a pixel from one position to another, with the price increasing with pixel distance.

Assuming $x \in \mathbb{R}^n_+$, are two non-negative data points in such a way that $\sum_i x_i = \sum_j y_j = 1$, inputs need to be normalized and let to the negative matrix a cost be $C \in \mathbb{R}^{n \times n}$ and $C_{i,j}$ be the cost of moving from $x_i$ to $y_j$, i.e., the pixel distance apart.

Therefore, the Wasserstein distance $d_w$ (x, y) can be defined.

$$d_w = (x,) = \min_{\Pi \in \mathbb{R}n \times n}(\Pi, C) \qquad \dots (1)$$

For $\Pi 1 = x,\ \Pi^T 1 = y$

Where $\Pi$ is the depreciation over a transport plan, whose inputs $\Pi_{i,j}$ shows how far mass moves from $x_i$ to $y_j$.

## 3.2    Projected Gradient Descent (PGD) attack

The Projected Gradient Descent attack is an iterative method in which, after each iteration, the perturbation is projected onto an lp-ball of a specified radius (in addition to clipping the values of the adversarial sample so that it lies in the permitted data range). This is the attack proposed by Madry et al. for adversarial training.

$$\mathcal{B}\ (x, \varepsilon) = \{z \in X|\ x - z\ |_p \leq\}\cdot \begin{cases} x' = x \\ x' = \Pi\ (\ x' = \psi(x', y)) \end{cases} \qquad \dots (2)$$

Where $\Pi_\varepsilon$ = projection on the ball $\mathcal{B}$ (x, $\varepsilon$), $\psi$ = perturbation function, y = target label yon wants to attack.

Gradient descent is a typical optimization algorithm for training machine learning models and neural networks. These models learn over time using training data and the cost function within gradient descent functions as a barometer, assessing its correctness with each iteration of parameter updates. Linear regression and logistic regression are two common examples of algorithms with coefficients that can be optimized using gradient descent.

## 3.3    Saliency Map Method attack

This method determines the spatial support for a specific class in each image. It is the earliest and most widely used explanation approach for evaluating convolutional neural network predictions.

### 3.3.1    How to Create Saliency Map

**1.** The basic properties of an image are extracted, such as color, orientation, and intensity from it.

2. To create a feature map, these processed images are used to develop Gaussian pyramids.

3. A saliency map is created by averaging all feature maps.

**The salience bias** refers to our inclination to pay greater attention to items or information that are more notable while dismissing those that are not.

### 3.3.2    Working of saliency maps method.

The gradients of the output over the input are used to create the saliency map. Saliency maps process images to differentiate visual features in images. For example, colored images are converted to black-and-white photos to analyze the richest colors.

**Static saliency** uses image characteristics and statistics to locate the image's regions of interest.

**Motion saliency**: Optical flow is used to detect motion in a video. Moving objects are deemed significant.

**Abjectness**: refers to the likelihood that an image window will cover an item. These techniques produce a collection of bounding boxes that indicate where an object in an image might be.

**Saliency analysis** determines how frequently any code appears, its importance, or both. High-value regulations increase understanding or are valuable in solving real-world issues. As a result, saliency analysis might reveal what is non-recurrent but potentially essential to a study's goals.

This approach uses the Jacobian matrix of outputs about inputs. Looking at this matrix, one can forecast how the output probabilities will behave if an input feature is significantly modified. Saliency maps were created to aid in visualizing the classification model prediction process. The map assigns a weight to each input feature x (i) (e.g., each pixel), causing the model to predict class c = y (x) = 1 major f (x) c 0), where f (x) is the victim model's SoftMax probability vector.

$$S^+ (x_{(i)}, c) = \begin{cases} 0 \; If \; \frac{\partial f(x)}{\partial x} < 0 \; or \; \sum_{c' \neq c} \frac{\partial f(x)}{\partial x} > 0 \\ -\frac{\partial f(x)}{\partial x} < 0 \cdot \sum_{c' \neq c} \frac{\partial f(x)}{\partial x} \quad Otherwise \end{cases} \quad \dots (3)$$

$S^+$ is the measure of how much x(i) positively correlates with c while also negatively correlates with all other classes $c^{'} = c$

$$S- (x(i), t) \; = 0 \; If \; \frac{\partial f(x)}{\partial x} > 0 \; or \; \sum_{c' \neq c} \frac{\partial f(x)}{\partial x} < 0$$

Both saliencies $S^-$ and $S^+$ are too demanding when applied to individual input features because the cumulative gradient contribution across all nontargeted classes $\sum_{ca \neq c} \frac{\partial f(x)(c')}{\partial x(i)}$ will frequently trigger the minimal saliency requirement.

**Results of the attacks on digits Keras model**

The table below shows the different attacks and their success rate on digits Keras

Table 1: Attack test

| Sample index | Input Label (conf) | Adversarial Label (conf) | % Euclid. dist | Adversarial Input | Success |
|---|---|---|---|---|---|

| 0 | Fraud (1.000) | benign (1.000) | 0.0 | [4462.00 -2.30 1.76 -0.36 2.33 …  0.04 -0.15 239.93] | True |
|---|---|---|---|---|---|

Table 2: Summary of the different model attacks.

| Attack names | Total runs | Successes (%) | Best score (attack_id) | Best parameters |
|---|---|---|---|---|
| saliency Map Method | 1 | 5 (100.0%) | 0.5 | Theta: 0.1, gamma: 1, batch_size: 1, etc. |
| Projected Gradient Descent Common | 1 | 5 (100.0%) | 0.6 | Targeted: false, Norm: infinity, esp_size: 0.1, etc. |
| Wasserstein | 1 | 5 (100.0%) | 1 | Regularization: 3000.0, kernel_size: 5, etc., |

**Labels** are things that we predict or that a trained machine is intended to provide. The label instructs users on what value should be typed in the corresponding input field.

**Adversarial label learning** trains the classifier to execute as expected. When given noisy and possibly associated labels

**Adversarial inputs** are inputs to machine learning models that an attacker has explicitly created to force the model to make a wrong decision. A negative example would be a corrupted version of a legitimate input contaminated by adding a minor disturbance.

**Discussion**

This section discusses attacks performed and the result obtained by different attacks on digits keras. The results, summarized in tables one and two, show that the attacks were successful, with the best score more significant than one. It also shows that these attacks are robust to adversarial defense and can fool the model and make it vulnerable to the attacker. Wasserstein attack has the highest best score of 1 and 100% success, as shown in table1 and aims to determine the perturbation that maximizes a model's loss on a given input. It is commonly stated as the perturbing's $L_2$ or $L_\infty$ norm, and it is added so that the content of the adversarial example is identical to the unperturbed sample.

For instance, Wasserstein's distance-bounded threat model limits disturbance to pixel-mass movements. The perturbations in the adversarial attack formed within a ball represent the image's actual content and structure. For example, in the top row, we can see a Wasserstein perturbation that hits all pixels indiscriminately rather than the space around. A standard model with binarization, a model provably resilient to perturbations of at most 0.1, and an adversarially trained model are all considered in MNIST and use only 100 iterations of projected gradient descent.

## Conclusion

Artificial intelligence (AI) and machine learning (ML) are two terms that have gained popularity in recent years from the software engineering standpoint. However, despite their widespread use, the community has paid little attention to AI and machine learning tools and applications. This paper presents cybersecurity issues in AI/ML tools. The attacks performed to show a success rate of 100%, with the best score of 0.5, 0.6, and 1 with different attacks.

**Reference**

[1] Atienza, Rowel. *Advanced Deep Learning with Keras: Apply deep learning techniques, autoencoders, GANs, variational autoencoders, deep reinforcement learning, policy gradients,* etc. Packt Publishing Ltd, 2018.

[2] Merino, Tim, et al. "Expansion of cyber-attack data from unbalanced datasets using generative adversarial networks." *International Conference on Software Engineering Research, Management, and Applications*. Springer, Cham, 2019.

[3] Li, Linyi, et al. "Sok: Certified robustness for deep neural networks." *arXiv preprint arXiv:2009.04131* (2020).

[4] Li, Jincheng, et al. "Internal Wasserstein distance for adversarial attack and defense." *arXiv preprint arXiv:2103.07598* (2021).

[5] Combey, Theo, et al. "Probabilistic Jacobian-Based Saliency Maps Attacks." *Machine Learning and Knowledge Extraction* 2.4 (2020): 558-578.

[6] Wiyatno, Rey, and Anqi Xu. "Maximal Jacobian-based saliency map attack." *arXiv preprint arXiv:1808.07945* (2018).

[7] Qureshi, Ayyaz Ul Haq, et al. "An adversarial approach for intrusion detection systems using Jacobian saliency map attacks (jsma) algorithm." *Computers* 9.3 (2020): 58.

[8] Wong, Eric, Frank Schmidt, and Zico Kolter. "Wasserstein adversarial examples via projected stinkhorn iterations." *International Conference on Machine Learning*. PMLR, 2019.

[9] Khedher, Hatem Ibn, Mohamed Ibn Khedher, and Makhlouf Hadji. "Mathematical Programming Approach for Adversarial Attack Modelling." *International Conference on Agents and Artificial Intelligence*. 2021.

[10] Manaswi, Navin Kumar. "Understanding and working with Keras." *Deep Learning with Applications Using Python*. Apress, Berkeley, CA, 2018. 31-43.

[11] Ketkar, Nikhil. "Introduction to Keras." *Deep learning with Python*. Apress, Berkeley, CA, 2017. 97-111.

[12] Vani, A. K., et al. "Using the Keras model for accurate and rapid gender identification through detection of facial features." *2020 Fourth International Conference on Computing Methodologies and Communication (ICCMC)*. IEEE, 2020.

[13] Ott, Jordan, et al. "A Fortran-Keras deep learning bridge for scientific computing." *Scientific Programming* 2020 (2020).

[14] Bhagwat, Ritesh, Mahla Abdolahnejad, and Matthew Moocarme. *Applied deep learning with Keras: Solve complex real-life problems with the simplicity of Keras*. Packt Publishing Ltd, 2019.

[15] Wang, Run, et al. "Amora: Black-box adversarial morphing attack." *Proceedings of the 28th ACM International Conference on Multimedia*. 2020.

[16] Atienza, Rowel. *Advanced Deep Learning with TensorFlow 2 and Keras: Apply DL, GANs, VAEs, deep RL, unsupervised learning, object detection, segmentation,* etc. Packt Publishing Ltd, 2020.

[17] Lanfredi, Ricardo Bigolin, Joyce D. Schroeder, and Tolga Tasdizen. "Quantifying the preferential direction of the model gradient in adversarial training with projected gradient descent." *arXiv preprint arXiv:2009.04709* (2020).

[18] Croce, Francesco, and Matthias Hein. "Sparse and imperceivable adversarial attacks." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019.

[19] Watney, M. M. "Artificial intelligence and its' legal risk to cybersecurity." *European conference on cyber warfare and security*. Academic Conferences International Limited, 2020.

[20] Chander, Bhanu, and Gopalakrishnan Kumaravelan. "Cyber security with AI—Part I." *The" Essence" of network security: An end-to-end panorama*. Springer, Singapore, 2021. 147-171.

[21] Ciolino, Matthew, Josh Kalin, and David Noever. "Fortify Machine Learning Production Systems: Detect and Classify Adversarial Attacks." *arXiv preprint arXiv:2102.09695* (2021).

**Appendices**

**Appendix A**: Test_attack.py

```
                                                           test_attacks_new.py
 Open    ▼   ⌐+                                          ~/counterfit/tests/counterfit/core

 1 import pytest
 2 from counterfit.core.attacks import CFAttack
 3 from counterfit.core.attacks import CFAttackOptions
 4
 5
 6 def test_cfattack_options():
 7     sample_options = {
 8         "option_int": 1,
 9         "option_float": 1.0,
10         "option_str": "hello"
11     }
12
13     cfattack_options = {
14         "sample_index": 0,
15         "logger": "default"
16     }
17
18     default_options_list = list(sample_options.keys())
19     cfattack_options_list = list(cfattack_options.keys())
20     all_options = {**sample_options, **cfattack_options}
21
22     attack_options = CFAttackOptions(
23         **all_options,
24         default_options_list=default_options_list,
25         cfattack_options_list=cfattack_options_list
26     )
27
28     assert attack_options is not None
29
30
31 def test_cfattack_get_default_options():
32     sample_options = {
33         "option_int": 1,
34         "option_float": 1.0,
35         "option_str": "hello"
36     }
37
38     cfattack_options = {
39         "sample_index": 0,
40         "logger": "default"
41     }
42
43     default_options_list = list(sample_options.keys())
44     cfattack_options_list = list(cfattack_options.keys())
45     all_options = {**sample_options, **cfattack_options}
46
47     attack_options = CFAttackOptions(
48         **all_options,
```

```
49            default_options_list=default_options_list,
50            cfattack_options_list=cfattack_options_list
51      )
52
53      default_options = attack_options.get_default_options()
54
55      assert sample_options == default_options
56
57
58  def test_cfattack_get_current_options():
59      sample_options = {
60          "option_int": 1,
61          "option_float": 1.0,
62          "option_str": "hello"
63      }
64
65      cfattack_options = {
66          "sample_index": 0,
67          "logger": "default"
68      }
69
70      default_options_list = list(sample_options.keys())
71      cfattack_options_list = list(cfattack_options.keys())
72      all_options = {**sample_options, **cfattack_options}
73
74      attack_options = CFAttackOptions(
75          **all_options,
76          default_options_list=default_options_list,
77          cfattack_options_list=cfattack_options_list
78      )
79
80      attack_options.set_options({"option_int": 5})
81      current_options = attack_options.get_current_options()
82
83      assert current_options["option_int"] == 5
84
```

```
85
86 def test_cfattack_get_all_options():
87     sample_options = {
88         "option_int": 1,
89         "option_float": 1.0,
90         "option_str": "hello"
91     }
92
93     cfattack_options = {
94         "sample_index": 0,
95         "logger": "default"
96     }
97
98     default_options_list = list(sample_options.keys())
99     cfattack_options_list = list(cfattack_options.keys())
100    all_options = {**sample_options, **cfattack_options}
101
102    attack_options = CFAttackOptions(
103        **all_options,
104        default_options_list=default_options_list,
105        cfattack_options_list=cfattack_options_list
106    )
107
108    all_options = attack_options.get_all_options()
109
110    assert all_options == {**sample_options, **cfattack_options}
111
112
113 def test_cfattack_save_previous_options():
114     sample_options = {
115         "option_int": 1,
116         "option_float": 1.0,
117         "option_str": "hello"
118     }
119
120     cfattack_options = {
121         "sample_index": 0,
122         "logger": "default"
123     }
124
125     default_options_list = list(sample_options.keys())
126     cfattack_options_list = list(cfattack_options.keys())
127     all_options = {**sample_options, **cfattack_options}
128
129     attack_options = CFAttackOptions(
130         **all_options,
```

**Appendix B**: Test_targets.py

```python
1 import pytest
2
3 from counterfit.core.targets import Target
4
5 from counterfit.targets.digits_blackbox.digits_blackbox import Digits
6 from counterfit.targets.digits_keras.digits_keras import DigitKeras
7 from counterfit.targets.creditfraud.creditfraud import CreditFraud
8 from counterfit.targets.satellite.satellite import SatelliteImagesTarget
9
10
11 @pytest.fixture(params=[CreditFraud, SatelliteImagesTarget])
12 def target(request):
13     yield request.param
14
15
16 def test_target_init(target):
17     new_target = target()
18     assert new_target is not None
19
20
21 def test_target_load(target):
22     new_target = target()
23     new_target.load()
24     new_target.set_loaded_status(status=True)
25     assert new_target.loaded_status == True
26
27
28 def test_target_get_samples(target):
29     new_target = target()
30     new_target.load()
31     new_target.set_loaded_status(status=True)
32     samples = new_target.get_samples(1)
33     assert samples is not None
34
35
36 def test_target_predict(target):
37     new_target = target()
38     new_target.load()
39     new_target.set_loaded_status(status=True)
40     samples = new_target.get_samples(1)
41     output = new_target.predict(samples)
42     assert output is not None
43
44
45 def test_target_outputs_to_labels(target):
46     new_target = target()
47     new_target.load()
48     new_target.set_loaded_status(status=True)
49     samples = new_target.get_samples(1)
50     output = new_target.predict(samples)
51     labels = new_target.outputs_to_labels(output)
52     assert labels is not None
```

## **Appendix C**: Digits_keras.py

```python
import os
import tensorflow as tf
from tensorflow import keras as K
import numpy as np

from counterfit.core.targets import Target

# for ART
tf.compat.v1.disable_eager_execution()


class DigitKeras(Target):
    target_data_type = "image"
    target_name = "digits_keras"
    target_endpoint = "mnist_model.h5"
    target_input_shape = (28, 28, 1)
    target_output_classes = ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"]
    target_classifier = "keras"
    X = []

    def load(self):
        if not os.path.isfile(self.fullpath(self.target_endpoint)):
            print("[!] Model file not found!")
            self.create_model()
        else:
            self.model = K.models.load_model(
                self.fullpath(self.target_endpoint))
        (train_x, train_y), (test_x, test_y) = K.datasets.mnist.load_data()

        self.X = test_x.astype(np.float32) / 255.  # float type [0,1]
        self.X = self.X.reshape(-1, 28, 28, 1)

    def create_model(self):
        # 0. get started
        print(
            f"   - Training new modeil. Begin MNIST using Keras {K.__version__}")
        np.random.seed(1)
        tf.random.set_seed(1)

        # 1. load data
        print("    - Loading train and test data ")
        (train_x, train_y), \
            (test_x, test_y) = K.datasets.mnist.load_data()
        train_x = train_x.reshape(60_000, 28, 28, 1)
        test_x = test_x.reshape(10_000, 28, 28, 1)
        train_x = train_x.astype(np.float32)
```

```python
47          test_x = test_x.astype(np.float32)
48          train_x /= 255
49          test_x /= 255
50          train_y = K.utils.to_categorical(train_y, 10)
51          test_y = K.utils.to_categorical(test_y, 10)
52
53          # self.X = [test_x]
54          self.test_x = test_x
55
56          # 2. define model
57          print(
58              "      - Creating network with two Convolution, two Dropout, two Dense layers ")
59          g_init = K.initializers.glorot_uniform(seed=1)
60          opt = K.optimizers.Adam(learning_rate=0.01)
61          x = K.layers.Input(shape=(28, 28, 1))
62          con1 = K.layers.Conv2D(
63              filters=32,
64              kernel_size=(3, 3),
65              kernel_initializer=g_init,
66              activation='relu',
67              padding='valid')(x)
68
69          con2 = K.layers.Conv2D(
70              filters=64,
71              kernel_size=(3, 3),
72              kernel_initializer=g_init,
73              activation='relu',
74              padding='valid')(con1)
75
76          mp1 = K.layers.MaxPooling2D(pool_size=(2, 2))(con2)
77          do1 = K.layers.Dropout(0.25)(mp1)
78          z = K.layers.Flatten()(do1)
79          fc1 = K.layers.Dense(
80              units=128,
81              kernel_initializer=g_init,
82              activation='relu')(z)
83
84          do2 = K.layers.Dropout(0.5)(fc1)
85          fc2 = K.layers.Dense(
86              units=10,
87              kernel_initializer=g_init,
88              activation='softmax')(do2)

90          model = K.models.Model(x, fc2)
91          model.compile(
92              loss='categorical_crossentropy',
93              optimizer=opt,
94              metrics=['accuracy'],
95              run_eagerly=False)
96
97          # 3. train model
98          batch_size = 100
99          max_epochs = 2
00          print(f"      - Starting training with batch size {batch_size}")
01          model.fit(
02              train_x,
03              train_y,
04              batch_size=batch_size,
05              epochs=max_epochs,
06              verbose=1)
07          print("      - Training finished ")
08
09          # 4. evaluate model
10          evaluation = model.evaluate(
11              test_x,
12              test_y,
13              verbose=0)
14
15          loss = evaluation[0]
16          acc = evaluation[1] * 100
17          print(f"      - Test data: loss = {loss}, accuracy = {acc}")
18
19          # 5. save model
20          print("      - Saving MNIST model to disk ")
21          model.save(self.fullpath("mnist_model.h5"))
22          self.model = model
23
24          print("[+] Done!")
25
26      def predict(self, x):
27          x = (x.reshape(-1, 28, 28, 1) * 255).astype(int).astype(float)/255.  # accept as quantized image
28          pred_probs = self.model.predict(x)
29          return pred_probs
```

**Appendix D**: Wasserstein attack



**Appendix E**: Projected Gradient Decent Common Attack

**Appendix F**: Saliency Map Method Attack