
Implementing a KNN based TTE attack detector

Tanmay Ambadkar

Department of Computer Science and Engineering
Pennsylvania State University
University Park, PA, 16803
ambadkar@psu.edu

Abstract

This report presents an implementation of a K*NN detector to detect TTE attacks. The detector first clusters the activations of the input data for every class present in the dataset and uses a KNN classifier on the cluster means obtained. This implementation yields a robust detector and can detect zeroth-order optimization (ZOO) attacked images with reasonable accuracy. Multiple experimental settings have been tested to select the best hyper-parameters and a decision statistic for the detector.

1 Introduction

The following assignment aims to implement a K*NN-based TTE attack detector. The model being used is the ResNet with four convolutional blocks and a linear layer for classification. The dataset used is the MNIST dataset, and the attack method is the ZOO attack. The report first explains the intuition behind using a K*NN detector. I explain the formulation of the detector next and the experimental setting to select the best hyper-parameters, which validate the intuition. The section after will explain the detection statistic, and the final section will show the results obtain.

The ZOO attack implementation and model used are from the last assignment.

2 Patch-wise ZOO-attack

In the last report, the patch-wise implementation was not clearly defined. In the proposed ZOO attack, I select a patch (3x3 or 4x4 block) and then optimise the image for that. Essentially, the coordinate descent is expanded to become a patch. What this does is increase the number of updates made to an image in a single optimization step, reducing the total steps required to attack an image. Another option would be to randomly select 9 coordinates in the image and optimize, which might yield the same results, but I have not experimented with the method. My attack methodology is illustrated in Figure 1

3 Intuition

When we train a deep learning model, the activations highlight the most important parts of the image, which can generate a good classification result. Every image belonging to a class will have similar characteristics, which can be identified by looking at the activations of the neural network. A nice visualization of activations has been provided by [1] (Figure 2)

At the shallower levels, we can make out the most important parts of the digit 2. The network captures the edges which define the shape of 2. As we go deeper, the representations are more abstract and cannot be understood by humans. These activation maps give us a hint on how to define our TTE detector.

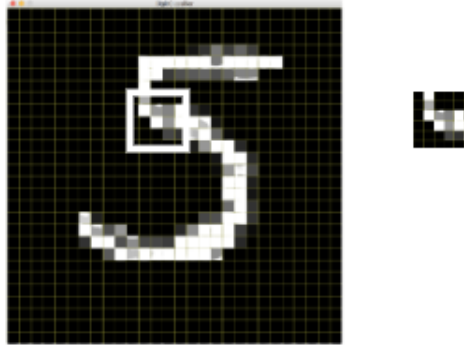


Figure 1: ZOO attack using Patch-wise gradient descent

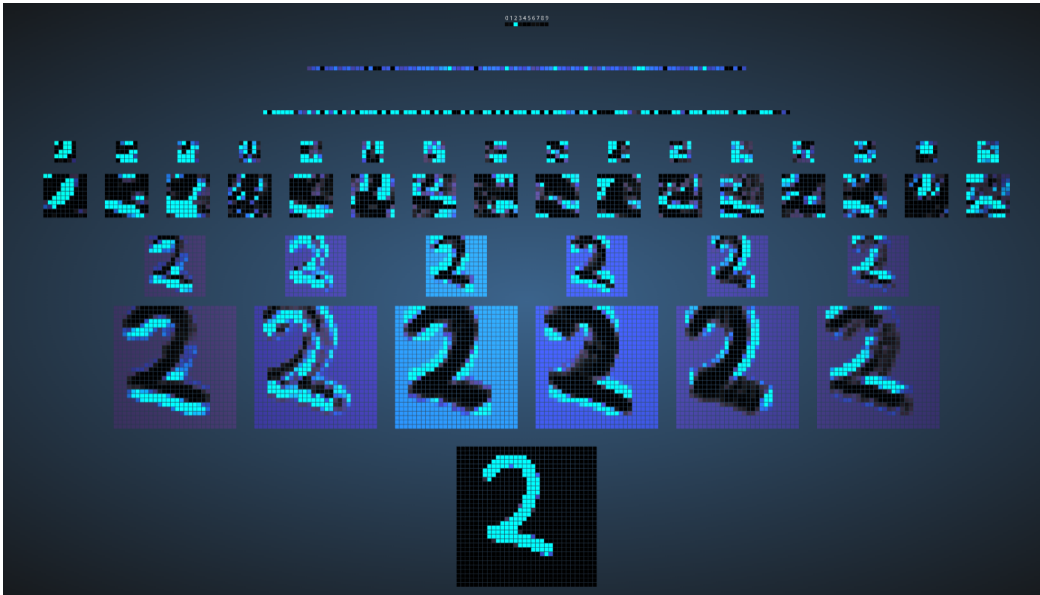


Figure 2: Activations of neural network for the MNIST digit 2

We can see that the activation maps for different digits will differ, so we can extract these activations and cluster them to get k different cluster centers for a class label $i \in 0, 9$. Then, any new activation will fall within one of these cluster centers. We can then use a distance-based voting rule to identify the class of any new image. This is the nearest neighbor rule popularly used in classification.

Upon inspection of the activations at different levels, we can infer that the activations at the shallower levels will have a lower distance in the cluster centers for different classes than at the deeper levels. This is because the maps differ less in the shallower layers than in the deeper ones, because the representations captured at deeper levels are more specific towards a digit, thus significantly different from each other. When a new point is to be classified, the distance between the cluster centers for the true class will be significantly smaller than the wrong class cluster centers for shallower layers than in deeper layers. Thus, it makes sense to use shallower layers to build the TTE detector than the deeper layers. Experimental results validate our hypothesis.

If we get an attacked image, the initial activations will still look like a non-attacked image. What happens at the deeper layers due to the attacked points causes a misclassification. Thus, when we use the nearest neighbor rule, the first neighbor will always return the correct class, while the next nearest neighbor might correspond to the attacked class.

In addition, we test the performance of our KNN classifier versus the DNN classifier on clean test sets to check the performance. Here, we will see that deeper layer activations give better performance. This is because at deeper layers, the embedded representations differ significantly to give a higher confidence, which in turn help the distance metric nearest neighbors classifier.

Thus, our intuition works in two ways. Deeper layers give better classification accuracy on clean sets than shallower layers by using fine-grained details, while shallower layers are more robust on attacked images by looking at coarse details.

I propose a detection statistic in the methodology section that uses the K*NN but discards the voting role, rather uses the distance and K neighbors to detect attacks with a better accuracy.

4 Methodology and TTE Detection statistic

Using the developed intuition, we first extract the activation maps for all convolutional blocks in the neural network. We have four such blocks in the network. We choose K clusters for every class i in the dataset and use the K-means algorithm to obtain cluster centers for the activations corresponding to the class i . Thus, we have a set of $K * N$ data points, where N is the number of classes. Using these data points, we create a nearest-neighbors classifier, which can classify non-attacked and attacked images with good accuracy, and thus becomes a robust classifier.

However, we can make the attack detector even better. For an image, we get K neighbors along with their distances. If the predicted label is present in one of the K neighbors, we find the one with the least distance and compare it with a threshold. If the distance is less than the threshold, the image is labeled as not attacked. If it is greater, it means that the image is attacked. This will mitigate false positives, where a digit might look like some other but not. For example, 6 and 0 can easily be misclassified.

To decide the threshold, we vary it from 0 to 100, plot a ROC curve, and measure the area under it. We find the point in the ROC curve closest to (0, 1) and use that as our decision threshold.

This goes against our intuition that shallower layers will perform better. In this method, the decision threshold and the distances calculated matter greatly, along with the number of neighbors being considered. The number of neighbors determines what points are to be considered. It is possible that the attack yields top 4 neighbors as the attacked class and the 5th neighbor as the true class. The KNN rule will fail here, giving us a false negative. Our rule will ensure that the 5th neighbor is considered and checked within a threshold. If the threshold check fails, then we can consider it as not attacked. If it falls within the threshold, then the image is attacked. Thus, the predicted label is required here.

The deeper layers will have clusters much farther from each other. Also, the ZOO attack implemented tries to misclassify by a very small margin. Thus, the activations in the initial layers will have all distances within a threshold and perform poorly. In the deeper layers, the distances will be greater. Thus, the threshold will greatly affect performance. Experimental results show that any layer will perform well, but the difference is not as large.

The steps to create the TTE detector are:

1. Use a validation set to get activations and their corresponding labels.
2. Perform K-means clustering on activations for a label i , and store cluster centers. Repeat this step for all classes in the dataset.
3. Using cluster centers and their corresponding labels, build a KNN classifier with k neighbors using the euclidean distance.
4. Test this classifier's performance on a test set and compare with DNN accuracy. Tune hyperparameters of K-means and KNN if necessary
5. Select a distance threshold between 0 and 100. This will be the TTE detection statistic.
6. For an image, get corresponding activations and pass through KNN to get top K neighbors. If the predicted DNN label is part of the top K neighbors, get the first occurring index and its distance.
7. If the distance is less than threshold, the image is not attacked. If $>$ than threshold, the image is attacked.

8. Vary this decision threshold and plot a ROC curve. The point where the $(fpr, tpr) \approx (0, 1)$ is the best threshold.

Thus, using this method, we have a robust classifier and an attack detector.

4.1 Why does this threshold method work?

One might question the logic or merit of using such a detection method. If the clusters are well-defined and separated by some hyperplane for each class, this method is not required. The K*NN rule will always work. Which is correct, except it is quite difficult to visualize embeddings in dimensions > 3 to prove the concept. We are creating clusters for 1 class rather than all classes; thus, cluster centers for 2 different classes will likely be very close. For the sake of simplicity, I will represent embeddings in 2 dimensions for 2 classes (Figure 3).

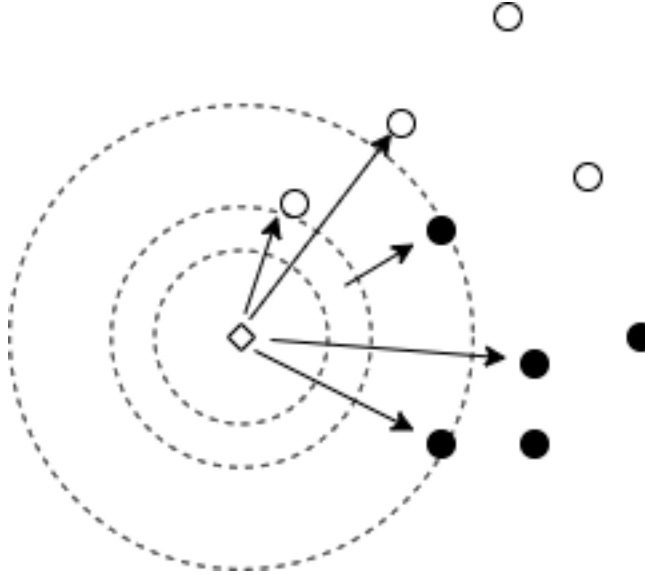


Figure 3: Illustration of detection statistic

The dashed circles represent our detection threshold. The black dots represent the target class of the attack, while the white dots represent the true class. The diamond represents the attacked image from the white class. If we choose a KNN based detector, we will get the same prediction as the DNN, and say that the image is not attacked. However, if we use the detection threshold, we see that the DNN predicted class is not present for the two inner circles. Thus, we can say that the image is attacked. We will conclude that the image is not attacked if we increase our tolerance (outermost circle). Thus, the threshold is an important hyperparameter for detecting attacks.

5 Experimental results

Since a validation set held out of training is not available directly, I select 1000 randomly sampled images from the test set to create the detector. The first set of results involves creating a KNN classifier and comparing its accuracy with the DNN classifier. These results include the robustness of the KNN classifier on attacked images, The final set of results shows the performance of the attack detector.

5.1 Classification results

We test multiple configurations of the KNN classifier by varying the layer activation, the number of clusters in K-means, and the number of neighbors in KNN. In total, we test 64 such configurations to get the best parameters. The trained DNN gives 98.5% accuracy. The results are part of figure 4:

		clean_accuracy	attacked_accuracy			clean_accuracy	attacked_accuracy
layer	n_clusters			layer	n_neighbors		
layer1	10	0.794776	0.737562	layer1	3	0.822139	0.742537
	15	0.810945	0.731343		5	0.796020	0.746269
	20	0.794776	0.728856		7	0.796020	0.722637
	25	0.804726	0.732587		9	0.791045	0.718905
layer2	10	0.840796	0.763682	layer2	3	0.893035	0.796020
	15	0.878109	0.783582		5	0.870647	0.792289
	20	0.889303	0.796020		7	0.871891	0.779851
	25	0.894279	0.808458		9	0.866915	0.783582
layer3	10	0.950249	0.819652	layer3	3	0.967662	0.855721
	15	0.950249	0.827114		5	0.952736	0.839552
	20	0.956468	0.839552		7	0.951493	0.823383
	25	0.963930	0.845771		9	0.949005	0.813433
layer4	10	0.980100	0.712687	layer4	3	0.982587	0.697761
	15	0.982587	0.743781		5	0.983831	0.716418
	20	0.983831	0.703980		7	0.982587	0.728856
	25	0.981343	0.707711		9	0.978856	0.725124

(a) Using clusters as groups

(b) using neighbors as groups

Figure 4: Results on clean and attacked datasets, grouped by the number of clusters and neighbors

We see here that the number of clusters or neighbors does not affect the accuracy of the clean or attacked sets. However, an important observation is that the clean set classification accuracy is highest for deeper layers, which supports our intuition. One interesting result is that layer 3 activations provide a more robust classification of attacked images. We could hypothesize that these activations encode the high level information best. If we use this classifier for detecting attacked images by using the following rule - if KNN label is not equal to the DNN label, then the image is attacked, otherwise not; we will get > 85% detection accuracy. This is because the DNN has 0% accuracy, so at least 85% of images will be detected as attacked. carry most high-level information, thus providing a more robust classification. Thus, I choose the number of clusters as 25 and the number of neighbors as 5, with layer 4 activations.

5.2 KNN-TTE detector results

For this experiment, we took 200 clean images and 200 attacked images to test the results of the KNN-TTE detector. To find the best detection threshold, we vary the threshold values from 0 to 100 and plot the ROC curve. The thresholds for which the (fpr, tpr) tuple is $\approx (0, 1)$ would be the best decision thresholds. The ROC curve can be seen in Figure 5. The AUC obtained is 1.0, meaning our detector can correctly identify all attacked images without false positives. The thresholds obtained are 7, 8, and 9. Beyond this, our true positives decrease, but false positives do not. When the threshold is lesser than 7, the true positives remain unchanged, but false positives increase.

We have also plotted the ROC curves for the first layer activations to see if the same works on shallower layers. The results do not differ by much, and the AUC score differs by a minute value. The same is for layer 3 activations. However, if we were to build an attack detection and mitigation strategy, I would select layer 3 activations, as the KNN can also correctly classify the attacked image. The thresholds stay unchanged, 7, 8 or 9.

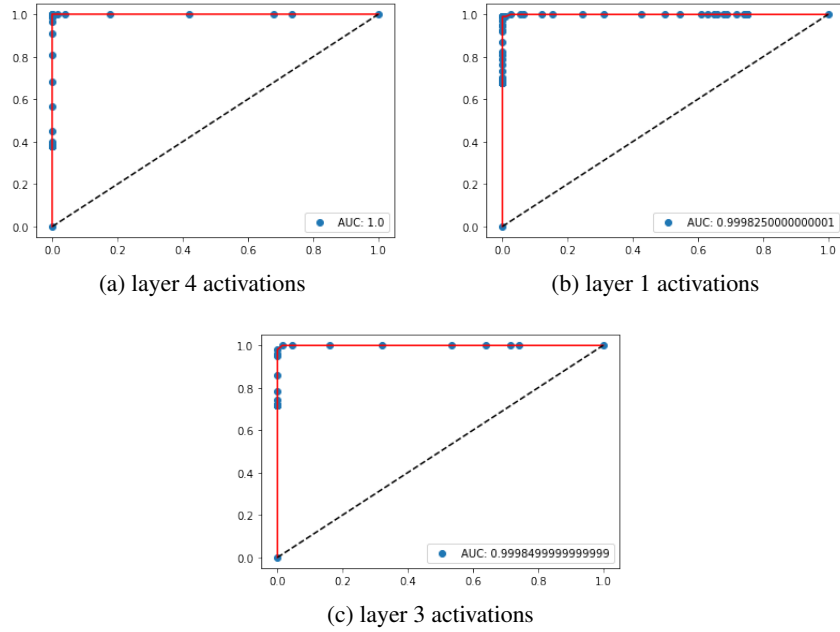


Figure 5: ROC curve for KNN TTE detector using various thresholds

6 Conclusion

The project proposes a KNN-based TTE attack detector. The KNN detector uses K-means cluster centers for activations to robustly classify attacked images. In practice, we achieve 85% accuracy on attacked images while the DNN achieves 0%. We use a threshold-based Nearest neighbor method to detect attacked images that achieves 0.998 AUC. The overall system can thus identify attacked images and robustly classify them.

Future work would be to identify why the layer 4 activations perform worse on attacked image classification using the KNN classifier than layer 3.

References

- [1] Adam W Harley. An interactive node-link visualization of convolutional neural networks. In *ISVC*, pages 867–877, 2015.

7 Appendix - Code

7.1 Model

The model code used to extract activations

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from CNN.resnet import ResNet18
5
6 class IntermediateRep(nn.Module):
7
8     def __init__(self, resnet, layer_to_extract):
9         super(IntermediateRep, self).__init__()
10        self.resnet = resnet
11        for param in self.resnet.parameters():
12            param.requires_grad = False
13        self.layer_to_extract = layer_to_extract
14
15    def forward(self, x):
16
17        resnet_out = torch.argmax(F.softmax(self.resnet(x), dim=1))
18        for name, layer in self.resnet.named_children():
19            x = layer(x)
20            if name == self.layer_to_extract:
21                break
22
23        return x, resnet_out
24
25 def create_model(layer = "layer1"):
26
27     model = ResNet18()
28     device = "cuda" if torch.cuda.is_available() else "cpu"
29     model.load_state_dict(torch.load('./model_weights/cpu_model.pth'))
30     model.to(device)
31     model.eval()
32
33     intermediate_rep = IntermediateRep(model, layer).cuda().eval()
34     intermediate_rep.eval()
35
36     return intermediate_rep
```

7.2 ZOO Attack

The following is the code used for the ZOO attack

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import numpy as np
5
6 device = "cuda" if torch.cuda.is_available() else "cpu"
7
8 def loss_fn(logits, label):
9
10    log_probs = F.log_softmax(logits, dim=1).reshape(-1,)
11    label_log_prob = log_probs[label].item()
12    log_probs[label] = -1e8
13    max_log_prob = log_probs.max()
14    return max(label_log_prob - max_log_prob, 0)
15
16
17 def calc_gradient(network, image, label, constant, h):
```

```

18
19     gradient = loss_fn(network(image+constant), label) - loss_fn(
20     network(image-constant), label)
21     gradient = gradient/h
22     return gradient
23
24 def zoo_attack(network, image, label):
25     '''
26
27     #todo you are required to complete this part
28     :param network: the model
29     :param image: one image with size: (1, 1, 32, 32) type: torch.
30     Tensor()
31     :param label: real label
32     :return: return a torch tensor (attack image) with size (1, 1, 32,
33     32)
34     '''
35     eta = 0.1
36     beta1 = 0.9
37     beta2 = 0.999
38     epsilon = 10**-8
39     alpha = 0.01
40     M = torch.zeros(*image.shape).to(device)
41     T = torch.zeros(*image.shape).to(device)
42     v = torch.zeros(*image.shape).to(device)
43     h = 0.001
44     f = loss_fn(network(image), label)
45     iter_ = 0
46     while f>0:
47         # for _ in range(1000):
48             iter_+=1
49             e = torch.zeros((1, 1, 32, 32))
50             r = np.random.randint(0, 31-3)
51             c = np.random.randint(0, 31-3)
52             e[0, 0, r:r+3, c:c+3] = h
53             e = e.to(device)
54             gradient = calc_gradient(network, image, label, e, h)
55
56             # pdb.set_trace()
57             # gradient = 0
58             T[0, 0, r:r+3, c:c+3] = T[0, 0, r:r+3, c:c+3]+1
59             M[0, 0, r:r+3, c:c+3] = beta1*M[0, 0, r:r+3, c:c+3] +
60             (1-beta1)*gradient
61             v[0, 0, r:r+3, c:c+3] = beta2*v[0, 0, r:r+3, c:c+3] +
62             (1-beta2)*gradient**2
63             M_hat = M[0, 0, r:r+3, c:c+3]/(1-beta1**T[0, 0, r:r+3, c:c+3])
64             v_hat = v[0, 0, r:r+3, c:c+3]/(1-beta2**T[0, 0, r:r+3, c:c+3])
65
66             image[0, 0, r:r+3, c:c+3] = image[0, 0, r:r+3, c:c+3] - eta*
67             M_hat/(v_hat**0.5 + epsilon)
68             f = loss_fn(network(image), label)=
69
70     return images

```

7.3 Creating KNN defense

The following code calculates the TPR and FPR

```

1 def true_positive_rate(true, pred):
2     return sum(pred[true == 1])/sum(true[true == 1])
3
4 def false_positive_rate(true, pred):
5     return sum(pred[true == 0])/len(true[true == 0])

```


The following code uses grid-search to find best parameters for defense. This finds accuracy on clean set and attacked set

```

1 collected_stats = {
2     "layer": [],
3     "n_clusters": [],
4     "n_neighbors": [],
5     "clean_accuracy": [],
6     "attacked_accuracy": [],
7     "dnn_attacked": 0
8 }
9 for layer in layers:
10
11     model = create_model(layer)
12     training_activations = []
13     training_labels = []
14
15     for i, data in enumerate(test_set):
16         if i == 1000:
17             break
18         image, label = data[0].reshape(1,*data[0].shape).to(device),
19         data[1]
20         training_activations.append(model(image)[0].reshape(1,-1).cpu()
21         ().numpy())
22         training_labels.append(data[1])
23
24     test_activations = []
25     test_labels = []
26     test_activations_attacked = []
27     dnn_attacked_labels = []
28
29     for i, data in tqdm(enumerate(test_set)):
30         if i < 1000:
31             continue
32         image, label = data[0].reshape(1,*data[0].shape).to(device),
33         data[1]
34         test_activations.append(model(image)[0].reshape(1,-1).cpu()
35         ().numpy())
36         adv_image = zoo_attack(network=model.resnet, image=image,
37         label=label)
38         acts, pred_label = model(adv_image)
39         test_activations_attacked.append(acts.reshape(1,-1).cpu()
40         ().numpy())
41         dnn_attacked_labels.append(pred_label.item())
42         test_labels.append(data[1])
43
44     if i == 1200:
45         break
46
47     training_labels = np.array(training_labels)
48     test_labels = np.array(test_labels)
49     training_activations = np.concatenate(training_activations, axis
50     =0)
51     test_activations = np.concatenate(test_activations, axis=0)
52     test_activations_attacked = np.concatenate(
53     test_activations_attacked, axis=0)
54     dnn_attacked_labels = np.array(dnn_attacked_labels)
55
56     for n_clusters in clusters:
57         for n_neighbors in neighbors:
58             knn_X = []
59             knn_y = []

```

```

55         for k_label in range(10):
56             kmeans = KMeans(n_clusters = n_clusters, random_state
= 0)
57             kmeans.fit(training_activations[training_labels ==
k_label])
58             knn_X.append(kmeans.cluster_centers_)
59             knn_y.append(np.array([k_label]*n_clusters))
60
61             knn_X = np.concatenate(knn_X, axis=0)
62             knn_y = np.concatenate(knn_y, axis=0)
63
64             knn_classifier = KNeighborsClassifier(n_neighbors =
n_neighbors)
65             knn_classifier.fit(knn_X, knn_y)
66
67             collected_stats["clean_accuracy"].append(accuracy_score(
test_labels, knn_classifier.predict(test_activations)))
68             collected_stats["attacked_accuracy"].append(accuracy_score
(test_labels, knn_classifier.predict(test_activations_attacked)))
69
70             collected_stats["layer"].append(layer)
71             collected_stats["n_clusters"].append(n_clusters)
72             collected_stats["n_neighbors"].append(n_neighbors)
73             collected_stats["dnn_attacked"] = accuracy_score(
test_labels, dnn_attacked_labels)

1
2 def return_knn(dataset, model, n_clusters, n_neighbors):
3
4     intermediate_reps, true_labels, = [], []
5     for i, value in enumerate(dataset):
6         if i == 500:
7             break
8         data_point = value[0].cuda().reshape(1, *value[0].shape)
9         activations, label = model(data_point)
10
11         if label!=value[1]:
12             continue
13         intermediate_reps.append(activations.reshape(-1,).cpu().numpy
())
14         true_labels.append(value[1])
15         del data_point
16
17         intermediate_reps = np.array(intermediate_reps)
18         true_labels = np.array(true_labels)
19
20         kmeans_models = []
21         for i in range(10):
22             kmeans = KMeans(n_clusters = n_clusters, random_state = 0)
23             kmeans.fit(intermediate_reps[true_labels == i])
24             kmeans_models.append(kmeans)
25
26
27         knn_X = []
28         knn_y = []
29
30         for i, kmeans in enumerate(kmeans_models):
31             knn_X.append(kmeans_models[i].cluster_centers_)
32             knn_y.append(np.array([i]*n_clusters))
33
34         knn_X = np.concatenate(knn_X, axis=0)
35         knn_y = np.concatenate(knn_y, axis=0)
36
37         knn_classifier = KNeighborsClassifier(n_neighbors = n_neighbors)
38         knn_classifier.fit(knn_X, knn_y)

```

```

39
40
41     return knn_classifier, knn_y

```

7.4 KNN TTE detection statistic

```

1 dist_pred_attacked = {i:[] for i in range(0,101)}
2 true_attacked = []
3 total = 0
4 for i, (images, labels) in tqdm(enumerate(testloader)):
5
6     images, labels = images.to(device), labels.to(device)
7     activations, pred_label = model(images)
8     if pred_label.item() != labels.item():
9         continue
10    total+=1
11    #Checking defense on unattacked images
12
13    true_attacked.append(0)
14    dist, label_idx, pred_label = knn_defense(images, model, defense)
15    least_dist = distance_label(dist[0], label_idx[0], pred_label.
16    item(), sample_labels)
17    for distance in range(0,101):
18        if least_dist < distance:
19            dist_pred_attacked[distance].append(0)
20        else:
21            dist_pred_attacked[distance].append(1)
22
23    #checking defense on attacked images
24
25    adv_image = zoo_attack(network=model.resnet, image=images, label=
26    labels)
27    adv_image = adv_image.to(device)
28    activations, pred_label = model(adv_image)
29
30    true_attacked.append(1)
31    dist, label_idx, pred_label = knn_defense(adv_image, model,
32    defense)
33    least_dist = distance_label(dist[0], label_idx[0], pred_label.
34    item(), sample_labels)
35    for distance in range(0,101):
36        if least_dist < distance:
37            dist_pred_attacked[distance].append(0)
38        else:
39            dist_pred_attacked[distance].append(1)
40
41    if total >= 200:
42        break

```