

IT581 Adversarial Machine Learning

Lab Assignment - 02

Group 05

2021 17013: Shaikh Faizan Ahmed

2021 17014: Sonam Bharti

Question 1.

Find the "Minimum L_2 Norm" adversarial attack to a two-class linear classifier by getting optimal x , solving:

$$\text{maximize}_x ||x - x_0||^2$$

subject to

$$w^T x + w_0 = 0$$

Answer

Using Lagrange multiplier of the constrained optimization is given by:

$$\mathcal{L}(x, \lambda) = \frac{1}{2} ||x - x_0||^2 + \lambda(w^T x + w_0). \quad (1)$$

The solution of the optimization is the saddle point (x^*, λ^*) such that $\nabla_x \mathcal{L} = 0$ and $\nabla_\lambda \mathcal{L} = 0$.

Taking derivative of equation (1) with respect to x and then to λ yields:

$$\nabla_x \mathcal{L} = x - x_0 + \lambda w = 0, \quad (2)$$

$$\nabla_\lambda \mathcal{L} = w^T x + w_0 = 0, \quad (3)$$

Multiplying the equation (2) by w^T yields :

$$w^T(x - x_0) + \lambda w^T w = 0$$

$$\Rightarrow w^T x - w^T x_0 + \lambda ||w||^2 = 0$$

$$\Rightarrow -w_0 - w^T x_0 + \lambda ||w||^2 = 0$$

.

Thus, the optimal λ is:

$$\lambda^* = \frac{(w^T x_0 + w_0)}{||w||^2}. \quad (4)$$

Correspondingly, the optimal x is:

$$x^* = x_0 - \lambda^* w = x_0 - \left(\frac{w^T x_0 + w_0}{||w||^2} \right) \frac{w}{||w||_2}. \quad (5)$$

Question 2.

The MNIST database of handwritten digits, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. It is a good database for people who want to try learning techniques and pattern recognition methods on real-world data while spending minimal efforts on pre-processing and formatting.

A. Perform Linear Discriminant Analysis (LDA) on the MNIST dataset* for binary classification and find weights and bias.

a. Plot confusion matrix and accuracy.

B. Perform Minimum l2 Norm adversarial attack to set of two digits (eg. 1 and 7, 3 and 8) and get updated test dataset and try to predict with the model used in A.

a. Plot confusion matrix and accuracy.

b. Plot 10 misclassified digits as image.

Comment your observation.

Answer

A.

code:

```
1 import numpy as np
2 import pandas as pd
3
4 import matplotlib.pyplot as plt
5 from pylab import *
6
7 from sklearn import metrics
8 from sklearn.metrics import accuracy_score
9 from sklearn.datasets import fetch_openml
10 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
11 from sklearn.model_selection import train_test_split
12
13 mnist = fetch_openml("mnist_784")
14
15 X_train, X_test, y_train, y_test = train_test_split(mnist.train, mnist.target, test_size=0.40,
16                                                    random_state=0)
17 lda = LDA(n_components=9)
18 X_train_r2 = lda.fit(X_train, y_train).transform(X_train)
19
20 print(X_train_r2)
21 y_pred = lda.predict(X_test)
22 print("Accuracy before attack: ", accuracy_score(y_test, y_pred))
23
24 print(
25     f"Classification report for classifier {lda}:\n"
26     f"{metrics.classification_report(y_test, y_pred)}\n"
27 )
28 w0 = lda.intercept_
29 w = lda.coef_
30 print(f"Bias: \n{w0}")
31 print(f"\n\nWeights: \n{w}")
32
33 disp = metrics.ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
34 disp.figure_.suptitle("Confusion Matrix")
35 print(f"Confusion matrix:\n{disp.confusion_matrix}")
36
37 plt.show()
```

Listing 1: Question 2(A)

Output:

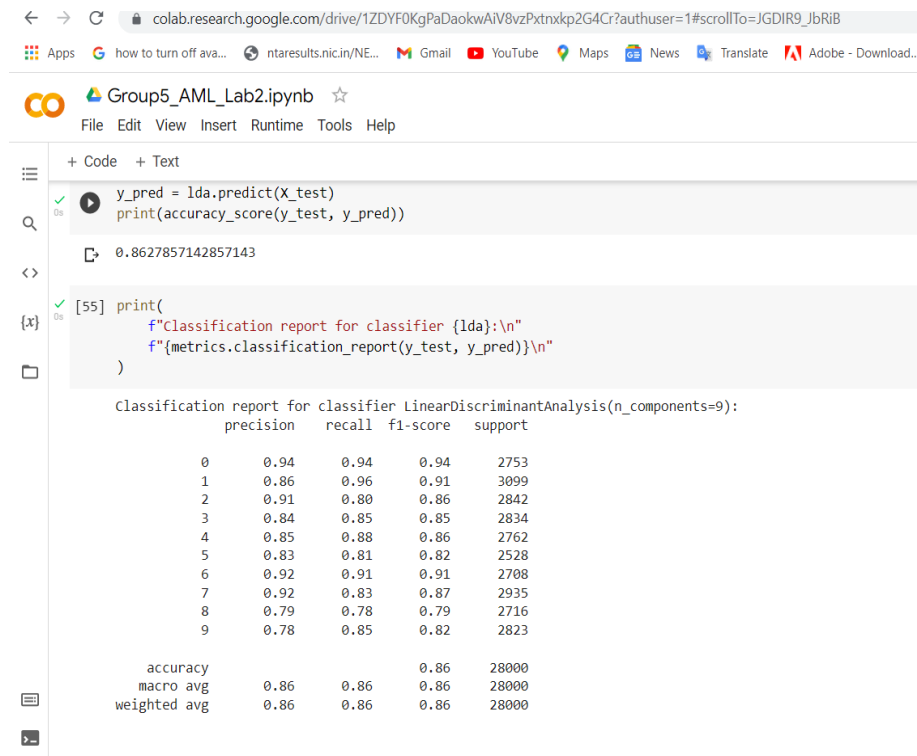


Figure 1: 2(A). Accuracy

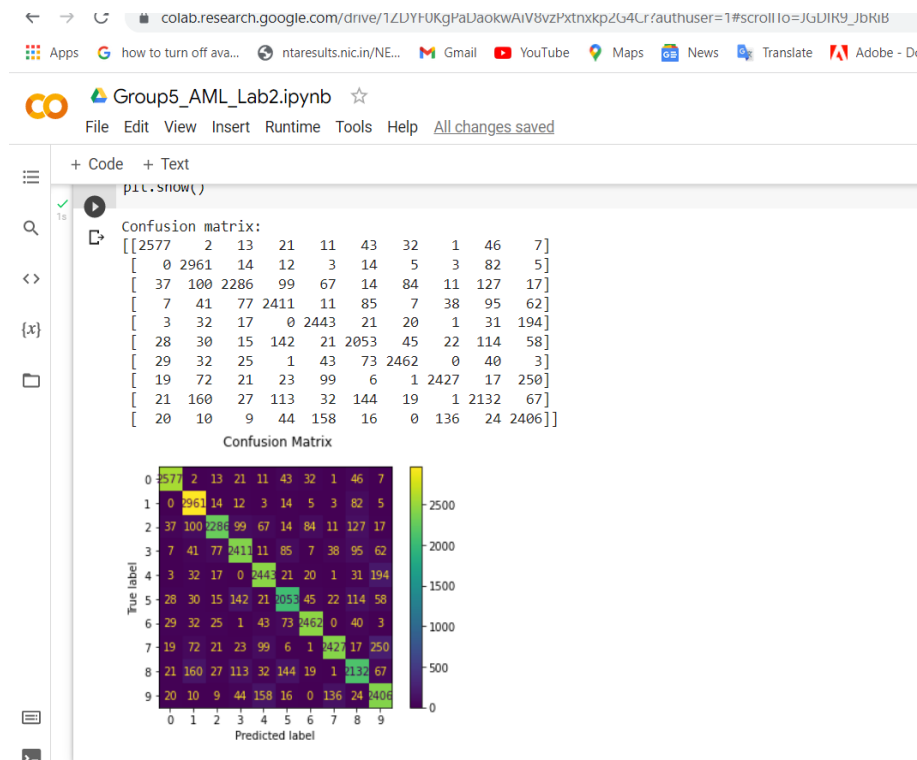


Figure 2: 2(A). Confusion Matrix

Definition:

A Confusion matrix is an $N \times N$ matrix used for evaluating the performance of a classification model, where N is the number of target classes. The matrix compares the actual target values with those predicted by the machine learning model. It is used to visualize important predictive analytics like recall, specificity, accuracy, and precision.

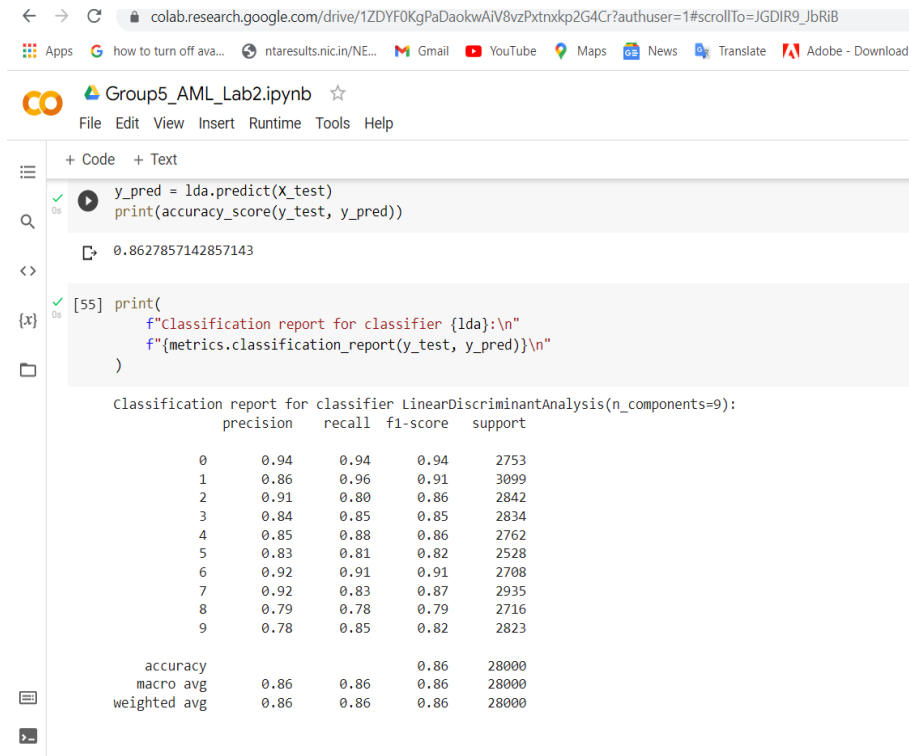


Figure 3: 2(A). Bias and Weights before attack

Bias is a phenomenon that skews the result of an algorithm in favor or against an idea. Bias is considered a systematic error that occurs in the machine learning model itself due to incorrect assumptions in the ML process. It is defined as the intercept of the function.

Whereas, Weights control the signal (or the strength of the connection) between two neurons. In other words, a weight decides how much influence the input will have on the output. It is defined as the coefficient.

Observation / Justification:

We have imported the **MNIST** dataset from "sklearn library". Then we split our train data and test data into 6:4 ratio as per asked. Then we applied LDA algorithm to train and testing the MNIST dataset. Then we calculated the accuracy of this algorithm by using `accuracy_score()` function and the accuracy rate over this dataset is approximately 87%.