

# 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  $(x^*, \lambda^*)$  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  $\lambda$  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  $\lambda$  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 import random
8
9 from sklearn import metrics
10 from sklearn.metrics import accuracy_score
11 from sklearn.model_selection import train_test_split
12 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
13
14 mnist_train = pd.read_csv("mnist_train.csv")
15
16 np.shape(mnist_train)
17
18 mnist_train.head()
19
20 mnist_train = mnist_train.loc[mnist_train['label'].isin([3, 8])]
21 X = mnist_train.iloc[:,1:785]
22 y = mnist_train.iloc[:,0]
23
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.40, random_state=0)
25
26 lda = LDA(n_components=1)
27 X_train_r2 = lda.fit(X_train, y_train)
28 y_pred = lda.predict(X_test)
29 print("Accuracy before attack: ",accuracy_score(y_test, y_pred))
30 print(y_pred.shape)
31
32 print(
33     f"Classification report for classifier {lda}:\n"
34     f"{metrics.classification_report(y_test, y_pred)}\n"
35 )
36
37 w0 = lda.intercept_
38 w = lda.coef_
39 print(f"Bias: \n{w0}")
40 print(f"\n\nWeights: \n{w}")
41
42 disp = metrics.ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
43 disp.figure_.suptitle("Confusion Matrix")
44 print(f"Confusion matrix:\n{disp.confusion_matrix}")
45
46 plt.show()
```

Listing 1: Question 2(A)

## Output:

```
print(y_pred.shape)

Accuracy before attack: 0.9524028512301679
(4349,)
```

```
print(
    f"Classification report for classifier {lda}:\n"
    f"{metrics.classification_report(y_test, y_pred)}\n"
)

w0 = lda.intercept_
w = lda.coef_
print(f"Bias: \n{w0}")
print(f"\n\nWeights: \n{w}")
```

Classification report for classifier LinearDiscriminantAnalysis(n\_components=1):

	precision	recall	f1-score	support
3	0.96	0.95	0.95	2213
8	0.95	0.95	0.95	2136
accuracy			0.95	4349
macro avg	0.95	0.95	0.95	4349
weighted avg	0.95	0.95	0.95	4349

Bias:  
[-3.28717619]

Figure 1: 2(A). Accuracy

```
Weights:
[[ 4.35017311e-15  1.93083418e-13  1.62663539e-13 -1.36722539e-13
 -4.87951163e-14  1.08176367e-13 -5.81587610e-14  2.25180615e-13
 -3.33698309e-14 -2.67936684e-15 -2.86147425e-14  2.11458786e-14
 6.03537382e-14 -1.94123527e-14  1.19153046e-14  2.17095997e-14
 4.02546064e-14 -5.46905533e-14  3.33369426e-14  7.66843667e-14
 5.16588103e-14  1.36755794e-14 -4.05060694e-14 -4.47083788e-14
 9.71080930e-14  1.27590666e-13  7.83462726e-15  2.05536551e-14
 1.17126285e-13  1.07776166e-13 -1.89727611e-15  1.76399323e-13
 5.64952361e-14  1.17477386e-13 -3.93019673e-14  4.13063286e-14
 9.13203574e-14 -9.60871798e-14  1.49566307e-13  3.66748645e-14
 -7.35065197e-15 -4.28326858e-14 -9.30200663e-14 -2.96740453e-14
 1.08773277e-13 -6.62196743e-14  3.02267321e-14  1.16475733e-13
 4.37513393e-14  1.01772886e-13  4.27909439e-15  1.27567402e-14
 4.08040789e-14  5.07127780e-14  1.33414467e-13 -2.26772437e-14
 -6.46156332e-14 -4.84097123e-14 -8.54604494e-14  5.39621645e-15
 1.21255291e-14 -1.14838823e-14  4.44665106e-14  1.31135286e-14
 -7.17097370e-14  1.61197843e-13 -1.22442025e-13 -1.76811395e-13
 1.04449163e-13 -8.95830160e-14 -3.40328161e-14 -9.94842326e-02
 -2.54718083e-02 -2.45176590e-02 -3.38217916e-03  2.69884662e-02
 -1.09589453e-13 -1.23518737e-13 -4.41829279e-15  4.86379037e-14
 4.18640526e-14 -4.04124690e-14 -1.35883026e-13  3.55970756e-14
 -1.78301528e-13  1.97792407e-15 -8.37600221e-14  8.52742652e-14
 2.27027555e-01 -1.61840444e-01  1.16308344e-01 -2.21831460e-02
 -3.49386642e-03  5.73639607e-03 -2.05098889e-03 -8.39384169e-03
 -5.46518998e-03  1.42011490e-02 -9.54329136e-03 -1.58150634e-02
 3.27235171e-03  4.28920748e-04  6.04811306e-03 -1.77066292e-02
 -1.73825266e-02 -1.32101570e-03  5.63313207e-02 -6.62056163e-02
 4.01487716e-01  7.32851656e-14 -1.02822742e-13  2.76555157e-14
 3.53511502e-14  8.26952313e-14  5.92947406e-02  1.41867203e-01
 9.65248436e-03 -7.53385862e-03  1.68660846e-03  1.16430208e-03]
```

Figure 2: 2(A). Weights

### 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.

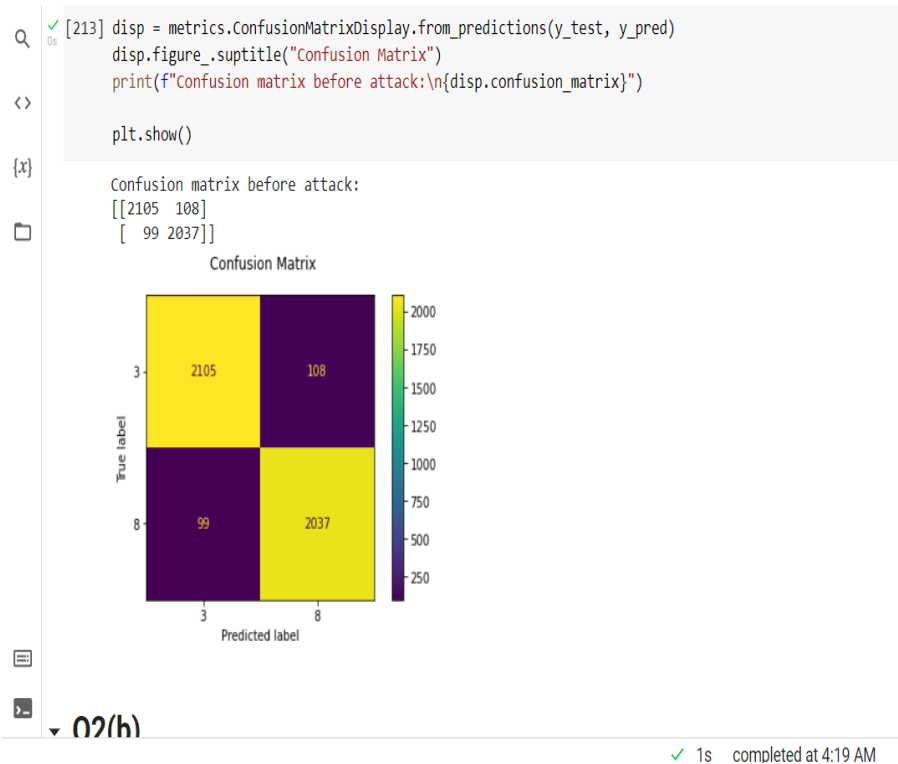


Figure 3: 2(A). Confusion Matrix before attack

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 95%.

B.a

code:

```
1 X_test_final = X_test.to_numpy(dtype = 'uint8')
2
3 X_attack=X_test_final-((X_test_final@np.transpose(w)+w0)@w)/np.linalg.norm(w)
4 print(X_attack.shape)
5
6 Y_attack=lda.predict(X_attack)
7 print(Y_attack.shape)
8
9 print("Accuracy after attack: ",accuracy_score(y_test , Y_attack))
10
11 print(f"Classification report for classifier {lda}:\n"
12       f"{metrics.classification_report(y_test, Y_attack)}\n")
13
14 disp = metrics.ConfusionMatrixDisplay.from_predictions(y_test, Y_attack)
15 disp.figure_.suptitle("Confusion Matrix")
16 print(f"Confusion matrix:\n{disp.confusion_matrix}")
17
18 plt.show()
```

Listing 2: Question 2(B).a

In this part we have to implement adversarial minimum norm attack over binary classifier. We tried to change 3 to 8 using minimum norm attack formula. And calculated accuracy score and Confusion Matrix after the attack.



Figure 4: 2(B).(a). Accuracy after attack

### Observation / Justification:

After the minimum  $L_2$  norm attack on the model, we observe the accuracy decreases to 4.7% for binary classifier. Also we can observe the changes in the output of confusion matrix.

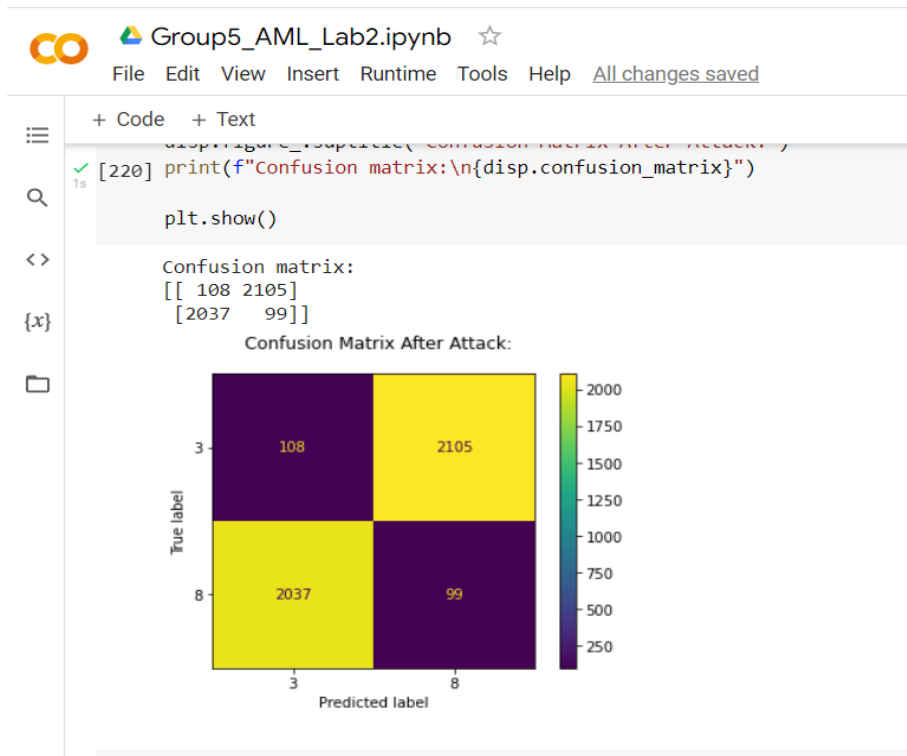


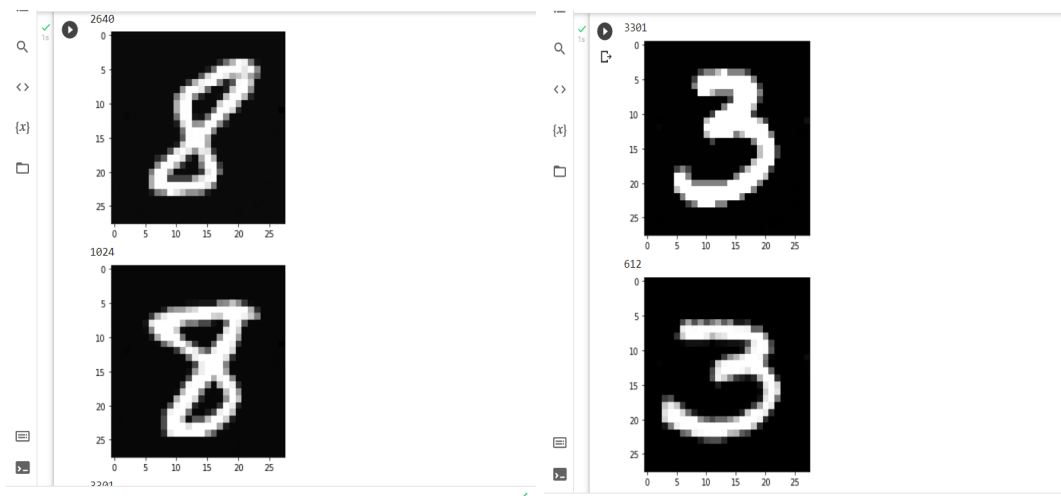
Figure 5: 2(B).(a). Confusion Matrix After attack

B.b

code:

```
1 X_attack_res=X_attack.reshape(X_attack.shape[0],28,28)
2 for i in range(0,10):
3     s=random.randint(0,X_attack.shape[0])
4     print(s)
5     plt.imshow(X_attack_res[s])
6     plt.show()
```

Listing 3: Question 2(B).b



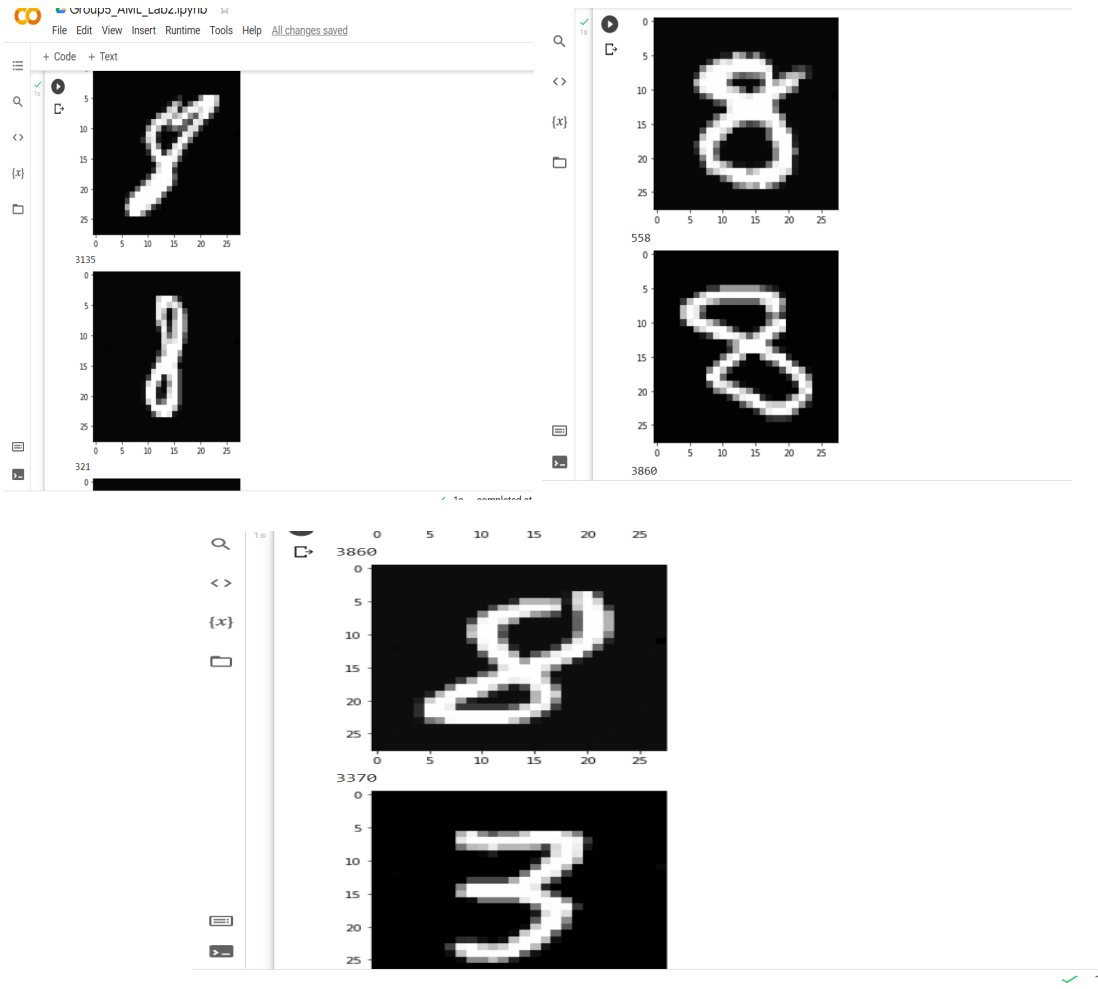


Figure 6: 2(B).(b). Image of misclassified digits

### Conclusion:

We can see that after minimum  $L_2$  norm attack most of the digit "3" is classified as "8" in 10 printed values.