

BANGLADESH UNIVERSITY OF ENGINEERING AND TECHNOLOGY

EEE 6002 Final Assignment

Selected Topics in Electrical and Electronic Engineering – Privacy Preserving
Machine Learning

Name: Mohammad Zunaed

ID: 0419062239

Email: 0419062239@eee.buet.ac.bd

Semester: April, 2020

1 Assignment Question

Final assignment instructions,

- Choose the MNIST dataset and consider digits 5 and 8 to be the two classes.
- Split the dataset into train, validation and test sets (e.g., using 70%, 10% and 20% splits). You can use sklearn library for splitting the data into these partitions.
- Write the Logistic Regression code by yourself (do not use any built-in classification function).
- Modify the Logistic Regression code such that it satisfies differential privacy (you can use either Laplace mechanism or Gaussian mechanism).
- Include the following plots for both non-private and private classifiers for $\epsilon = 0.01$ and $\delta = 1e - 5$:
 - Training and validation loss vs iterations
 - Norm of the gradient vs iterations
 - Training and validation loss vs number of training samples
 - Training and validation loss vs regularization coefficient λ
 - Required number of iterations vs step size α
 - Precision and Recall vs threshold
 - True Positive Rate vs False Positive Rate (ROC curve)
- Plot the classification accuracy against various ϵ values for fixed δ and sample size.
- Plot the classification accuracy against varying sample size for a fixed ϵ and δ .
- The report must contain a brief description of the dataset and the complete code.

2 Dataset Description

The MNIST dataset is a large dataset of handwritten digits widely used in the field of machine learning. The images in the dataset were normalized to fit into a 28x28 pixel bounding box and anti-aliased. Their color map is grayscale. The MNIST database contains 60,000 training images and 10,000 testing images. We downloaded the dataset from the PyTorch repository using their built-in function and merge the files. The MNIST dataset contains ten digits. However, for our problem, we need only the digits 5 and 8. For the gaussian mechanism, we normalize the dataset by 2-norm, and for the laplacian mechanism, we normalize it by 1-norm. Afterward, we split the dataset into train, validation, and test dataset according to the portion given in the assignment. The data distribution is shown in the table 1. Each of the data has 785 (including bias term) features. A few examples are shown in the figure 1.

Table 1: DATA DISTRIBUTION.

Type	Train	Validation	Test
Sample Number	9459	1051	2628
Percentage	72	8	20

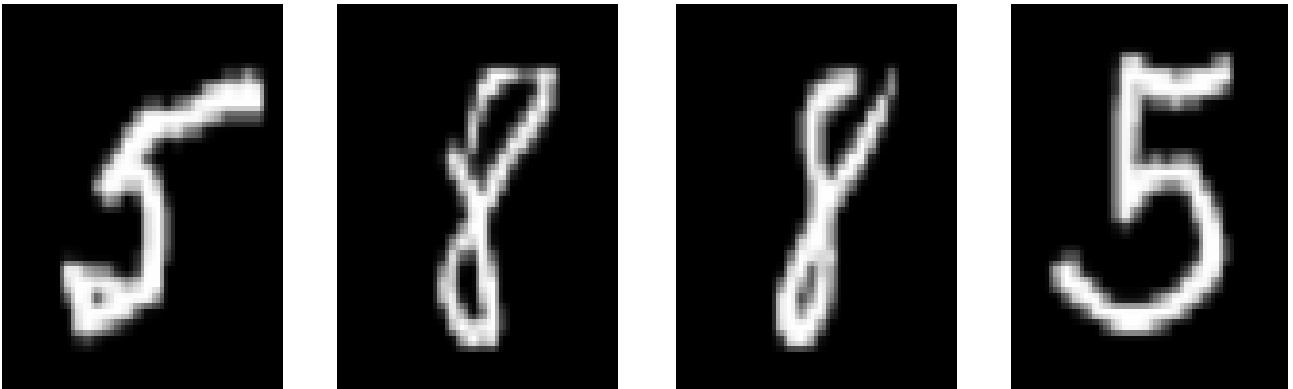


Figure 1: MNIST DATASET EXAMPLES.

Required libraries for my implementation are:

```
1 import os
2 import time
3 import torch
4 import random
5 import numpy as np
6 from tqdm import tqdm
7 import matplotlib.pyplot as plt
8 from torchvision import datasets, transforms
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import precision_score, recall_score, roc_curve,
    roc_auc_score, accuracy_score
```

The function that I use to process the dataset is given next:

```
1 def get_data(do_norm=False, norm_order=2):
2     transform = transforms.Compose([transforms.ToTensor()])
3     train_data = datasets.MNIST(root='data', train=True, download=True,
4     transform=transform)
5     test_data = datasets.MNIST(root='data', train=False, download=True,
6     transform=transform)
7
8     train_loader = torch.utils.data.DataLoader(train_data, batch_size=512)
9     test_loader = torch.utils.data.DataLoader(test_data, batch_size=512)
10
11     X_all = []; y_all = [];
12
13     for _, (X,y) in enumerate(train_loader):
14         X_all.append(X.squeeze(1).view(X.size(0), -1).data.numpy())
15         y_all.append(y.data.numpy())
16
17     for _, (X,y) in enumerate(test_loader):
18         X_all.append(X.squeeze(1).view(X.size(0), -1).data.numpy())
19         y_all.append(y.data.numpy())
20
21     X_all = np.concatenate(X_all)
22     y_all = np.concatenate(y_all)
23
24     X_all = X_all[np.logical_or(y_all == 5, y_all == 8)]
25     y_all = y_all[np.logical_or(y_all == 5, y_all == 8)]
26     y_all[y_all == 5] = 0
27     y_all[y_all == 8] = 1
28
29     if do_norm:
30         X_all_norm = np.linalg.norm(X_all, norm_order, axis=-1)
31         X_all_norm_max = np.max(X_all_norm)
32         X_all = X_all / X_all_norm_max
33
34     X, X_test, y, y_test = train_test_split(X_all, y_all, test_size=0.20,
35     random_state=42)
36     X_train, X_validation, y_train, y_validation = train_test_split(X, y,
37     test_size=0.10, random_state=42)
38
39     return X_train, y_train, X_validation, y_validation, X_test, y_test
```

3 Logistic Regression with Differential Privacy mechanism

The code for Logistic Regression has been implemented from scratch. The code files are uploaded in the following github repository https://github.com/rafizunaed/BUET-EEE_6002-PPML-Assignments.

Suppose the coefficients are ω , input features are X , ground truth is y , regularization coefficient is λ , m is the total number of data, n is the number of coefficients, r is the feature dimension, and α is the step size. The algorithm for logistic regression is:

$$\text{logits}, f(x) = \omega X = \omega_1 x_0 + \omega_2 x_1 + \omega_3 x_2 + \omega_4 x_3 + \dots + \omega_n x_r \quad (1)$$

$$\text{sigmoid function}, \sigma(z) = \frac{1}{1 + \exp(-z)} \quad (2)$$

$$\text{hypothesis}, \hat{y}(\omega, X) = \sigma(\omega X) \quad (3)$$

$$\text{loss}, L(y, \hat{y}, \omega) = - \left[\frac{1}{m} \sum_{i=1}^m (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) - \frac{\lambda}{2m} \sum_{i=1}^n \|\omega_i\|_2^2 \right] \quad (4)$$

$$\text{gradient}, \nabla_{\omega} L = \frac{1}{m} (X^T (\hat{y} - y) + \lambda \omega) \quad (5)$$

$$\text{weight update}, \omega = \omega - \alpha \nabla_{\omega} L \quad (6)$$

Here, ω_1 is the bias term and we will assume x_0 is 1. In the case of the Gaussian or Laplacian differential privacy mechanism, the only change will occur in the gradient calculation phase. We will add a noise term with the gradient before updating the weight values. For the Gaussian mechanism, the data will be norm-2 normalized. The updated steps for the Gaussian mechanism are:

$$\text{sensitivity}, \Delta_f = \frac{2}{m} \quad (7)$$

$$\text{gaussian noise}, n_g = N(0, \frac{2\Delta_f^2}{\epsilon^2} \log(1.25/\delta)) \quad (8)$$

$$\text{new gradient}, (\nabla_{\omega} L)_{new} = \nabla_{\omega} L + n_g \quad (9)$$

$$\text{weight update}, \omega = \omega - \alpha (\nabla_{\omega} L)_{new} \quad (10)$$

Here, N is normal distribution with 0 mean value. For the Laplacian mechanism, the data will be norm-1 normalized. The updated steps for the Laplacian mechanism are:

$$\text{sensitivity}, \Delta_f = \frac{2}{m} \quad (11)$$

$$\text{beta}, \beta = \frac{\Delta_f}{\epsilon} \quad (12)$$

$$\text{laplacian noise}, n_l = L_p(0, \beta) \quad (13)$$

$$\text{new gradient}, (\nabla_{\omega} L)_{new} = \nabla_{\omega} L + n_l \quad (14)$$

$$\text{weight update}, \omega = \omega - \alpha (\nabla_{\omega} L)_{new} \quad (15)$$

Here, L_p is laplacian distribution with 0 mean value.

The functions implementing the above algorithms are given below:

```
1 def sigmoid(x):
2     return 1 / (1 + np.exp(-x))
3
4 def calculate_loss(features, target, coefficients,
5     regularization_coefficient):
6     logits = np.dot(features, coefficients)
7     logits_sigmoid = sigmoid(logits)
8     loss = (- np.sum(target*np.log(logits_sigmoid)+ (1-target)*np.log(1-
9     logits_sigmoid)) \
10         + 0.5*regularization_coefficient*np.dot(coefficients,
11     coefficients))/len(target)
12     return loss
13
14 def train_one_iteration_dp_laplace(features, target, coefficients,
15     learning_rate,
16         regularization_coefficient, epsilon
17     =0.01):
18     f_x = np.dot(features, coefficients)
19     p_x = sigmoid(f_x)
20     gradient = (np.dot(features.T, (p_x - target)) +
21     regularization_coefficient * coefficients)/features.shape[0]
22     sensitivity = 2 / features.shape[0]
23     beta = sensitivity / epsilon
24     laplace_noise = np.random.laplace(0, beta, gradient.shape[0])
25     gradient_new = gradient + laplace_noise
26     coefficients -= learning_rate * gradient_new
27     return coefficients, gradient
28
29 def train_one_iteration_dp_gaussian(features, target, coefficients,
30     learning_rate,
31         regularization_coefficient, epsilon
32     =0.01, delta=1e-5):
33     sigma = (2/features.shape[0]) * np.sqrt(2*np.log(1.25/delta)) / epsilon
34     gaussian_noise = np.random.normal(loc=0.0, scale=sigma*sigma, size=
35     coefficients.shape[0])
36     f_x = np.dot(features, coefficients)
37     p_x = sigmoid(f_x)
38     gradient = (np.dot(features.T, (p_x - target)) +
39     regularization_coefficient * coefficients)/features.shape[0]
40     gradient_new = gradient + gaussian_noise
41     coefficients -= learning_rate * gradient_new
42     return coefficients, gradient
43
44 def train_one_iteration_non_dp(features, target, coefficients, learning_rate
45     , regularization_coefficient):
46     f_x = np.dot(features, coefficients)
47     p_x = sigmoid(f_x)
48     gradient = (np.dot(features.T, (p_x - target)) +
49     regularization_coefficient * coefficients)/features.shape[0]
50     coefficients -= learning_rate * gradient
51     return coefficients, gradient
```

```

1 def logistic_regression(features_train, target_train, features_validation,
2   target_validation, add_x0,
3   iterations, learning_rate,
4   regularization_coefficient, method='non_dp', epsilon=0.01, delta=1e-5):
5   set_random_state(42)
6
7   if add_x0:
8       x0 = np.ones((features_train.shape[0], 1))
9       features_train = np.hstack((x0, features_train))
10      x0 = np.ones((features_validation.shape[0], 1))
11      features_validation = np.hstack((x0, features_validation))
12      coefficients = np.zeros(features_train.shape[1])
13
14      best_coefficients = np.zeros(features_train.shape[1])
15      best_loss = 9999
16      count = 0
17      patience = 100
18      best_iteration = 0
19
20      all_train_loss=[]
21      all_validation_loss=[]
22      all_gradients=[]
23      progress_bar = tqdm(total=iterations, unit='iterations')
24      for step in range(iterations):
25          if method == 'non_dp':
26              coefficients, gradient = train_one_iteration_non_dp(
27                  features_train, target_train, coefficients, learning_rate,
28                  regularization_coefficient)
29          elif method == 'gaussian':
30              coefficients, gradient = train_one_iteration_dp_gaussian(
31                  features_train, target_train, coefficients, learning_rate,
32                  regularization_coefficient, epsilon, delta)
33          elif method == 'laplace':
34              coefficients, gradient = train_one_iteration_dp_laplace(
35                  features_train, target_train, coefficients, learning_rate,
36                  regularization_coefficient, epsilon)
37
38          train_loss = calculate_loss(features_train, target_train,
39              coefficients, regularization_coefficient)
40          validation_loss = calculate_loss(features_validation,
41              target_validation, coefficients, regularization_coefficient)
42
43          if validation_loss < best_loss:
44              best_loss = validation_loss
45              best_coefficients = coefficients.copy()
46              count = 0
47              best_iteration = step
48          else:
49              count += 1
50
51          progress_bar.set_postfix({'Train loss':train_loss, ' Validation_loss'
52              :validation_loss})
53          progress_bar.update(1)
54
55          all_train_loss.append(train_loss)
56          all_validation_loss.append(validation_loss)
57          all_gradients.append(gradient)
58
59          if count == patience:
60              learning_rate = learning_rate * 0.5

```

```

50         count = 0
51
52     progress_bar.close()
53
54     all_gradients = np.stack(all_gradients)
55     all_gradients_norm = np.linalg.norm(all_gradients, axis=-1)
56
57     return best_coefficients, all_train_loss, all_validation_loss,
all_gradients_norm, best_iteration

```

4 Plots

First, we will plot loss vs iterations, the norm of the gradients, precision and recall scores vs threshold, and the ROC curve for both non-differential and differential privacy methods. We calculate and plot these figures by running the following function:

```

1 def train(add_x0=True, iterations=5000, learning_rate=0.5,
2         regularization_coefficient=0.001, method='non_dp',
3         epsilon=0.01, delta=1e-5, data_norm=False, norm_order=2):
4
5     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data(
6         data_norm, norm_order)
7     best_coefficients, all_train_loss, all_validation_loss,
8     all_gradients_norm, _ = logistic_regression(
9         X_train, y_train,
10        X_validation, y_validation,
11        add_x0,
12        iterations,
13        learning_rate,
14        regularization_coefficient,
15        method,
16        epsilon,
17        delta,
18    )
19
20     os.makedirs('./logs/', exist_ok=True)
21     # plot of 5a. train and validation loss vs iterations
22     iterations = np.arange(0, iterations)
23     plt.plot(iterations, all_train_loss, label='train_loss')
24     plt.plot(iterations, all_validation_loss, label='validation_loss')
25     plt.legend()
26     plt.xlabel('iterations')
27     plt.ylabel('loss')
28     plt.title('loss vs iterations [5a]')
29     plt.savefig("./logs/{}_5a_loss_vs_iterations.png".format(method),
30         bbox_inches='tight', pad_inches=0)
31     plt.show()
32
33     # plot of 5b. Norm of the gradient vs iterations
34     plt.plot(iterations, all_gradients_norm)
35     plt.xlabel('iterations')
36     plt.ylabel('Norm of the gradient')
37     plt.title('Norm of the gradient vs iterations [5b]')
38     plt.savefig("./logs/{}_5b_Norm_of_the_gradient_vs_iterations.png".format(
39         method), bbox_inches='tight', pad_inches=0)
40     plt.show()
41
42     if add_x0:
43         X_test = np.hstack((np.ones((X_test.shape[0], 1)), X_test))

```



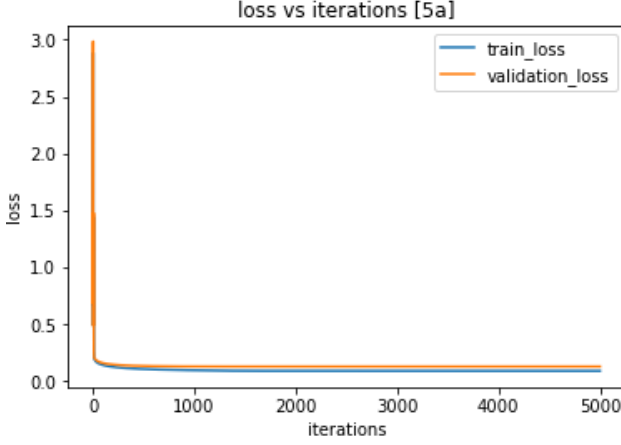
```

39 logits = np.dot(X_test, best_coefficients)
40 prediction_probabilites = sigmoid(logits)
41
42 all_precisions=[]
43 all_recalls=[]
44 all_thresholds=np.arange(0, 1, 0.01)
45 for threshold in all_thresholds:
46     y_true = y_test.copy()
47     y_pred = prediction_probabilites >= threshold
48     all_precisions.append(precision_score(y_true, y_pred, average='macro
'))
49     all_recalls.append(recall_score(y_true, y_pred, average='macro'))
50
51 # plot of 5f. Precision and Recall vs threshold
52 plt.plot(all_thresholds, all_precisions, label='Precision')
53 plt.plot(all_thresholds, all_recalls, label='Recall')
54 plt.legend()
55 plt.xlabel('Thresholds')
56 plt.ylabel('Precision and Recall Score')
57 plt.title('Precision and Recall Score vs Thresholds [5f]')
58 plt.savefig("./logs/{}_5f_Precision_and_Recall_Score_vs_Thresholds.png".
format(method), bbox_inches='tight', pad_inches=0)
59 plt.show()
60
61 # plot of 5g. True Positive Rate vs False Positive Rate (ROC curve)
62 fpr, tpr, _ = roc_curve(y_true, prediction_probabilites)
63 roc_score = roc_auc_score(y_true, prediction_probabilites, 'macro')
64 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area =
 {:.3f})'.format(roc_score))
65 plt.legend(loc="lower right")
66 plt.xlabel('False Positive Rate')
67 plt.ylabel('True Positive Rate')
68 plt.title('True Positive Rate vs False Positive Rate (ROC curve) [5g]')
69 plt.savefig("./logs/{}_5g_True_Positive_Rate_vs_False_Positive_Rate.png"
.format(method), bbox_inches='tight', pad_inches=0)
70 plt.show()
71
72 log_dict = {
73     'all_train_loss': all_train_loss,
74     'all_validation_loss': all_validation_loss,
75     'all_gradients_norm': all_gradients_norm,
76     'iterations': iterations,
77     'all_thresholds': all_thresholds,
78     'all_precisions': all_precisions,
79     'all_recalls': all_recalls,
80     'fpr': fpr,
81     'tpr': tpr,
82     'roc_score': roc_score,
83 }
84 np.save('./logs/{}_log_dict_1.npy'.format(method), log_dict)
85 if __name__ == '__main__':
86     # fig a,b,f,g for non_dp, gaussian, laplace
87     train(iterations=5000, method='non_dp', data_norm=False)
88     train(iterations=5000, method='gaussian', data_norm=True, norm_order=2)
89     train(iterations=5000, learning_rate=5, regularization_coefficient
=0.001, method='laplace', data_norm=True, norm_order=1)

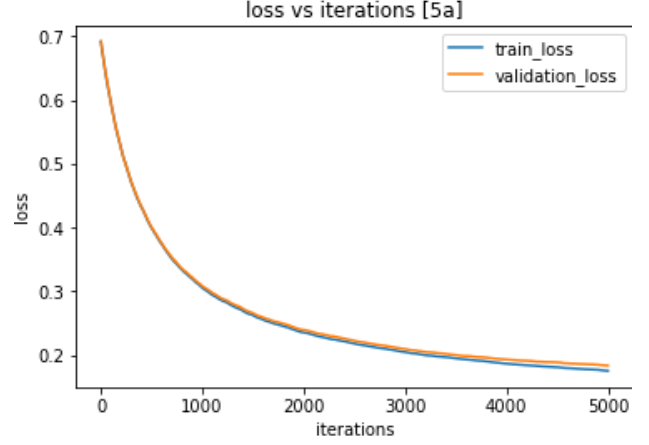
```

4.1 Training and Validation loss vs iterations

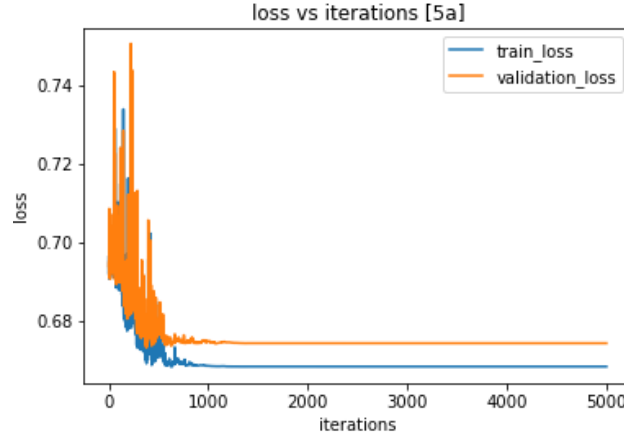
From the figure 2, we can see that the non-differential privacy method quickly converges but the gaussian method takes more iterations to converge. Also, the non-differential privacy method has a lower loss value than the gaussian method. However, the laplacian method has a very high loss value and also does not converge smoothly. In all the cases, validation loss is slightly higher than the training losses which indicates slight overfitting.



(a) Without Differential Privacy Method



(b) Gaussian Differential Privacy Method



(c) Laplacian Differential Privacy Method

Figure 2: Training and Validation loss vs iterations.

4.2 Norm of the gradient vs iterations

Now, we plot the Norm of the gradient against iterations in the figure 3. We see that, with the increasing number of iterations, the norm of the gradient is decreasing. This means that our model is converging to minimize the loss. For the non-differential privacy method, the curve is smooth. But for the differential methods, curves are not smooth due to added noise. The Gaussian mechanism is much more stable and close to the non-differential method than the Laplacian mechanism.

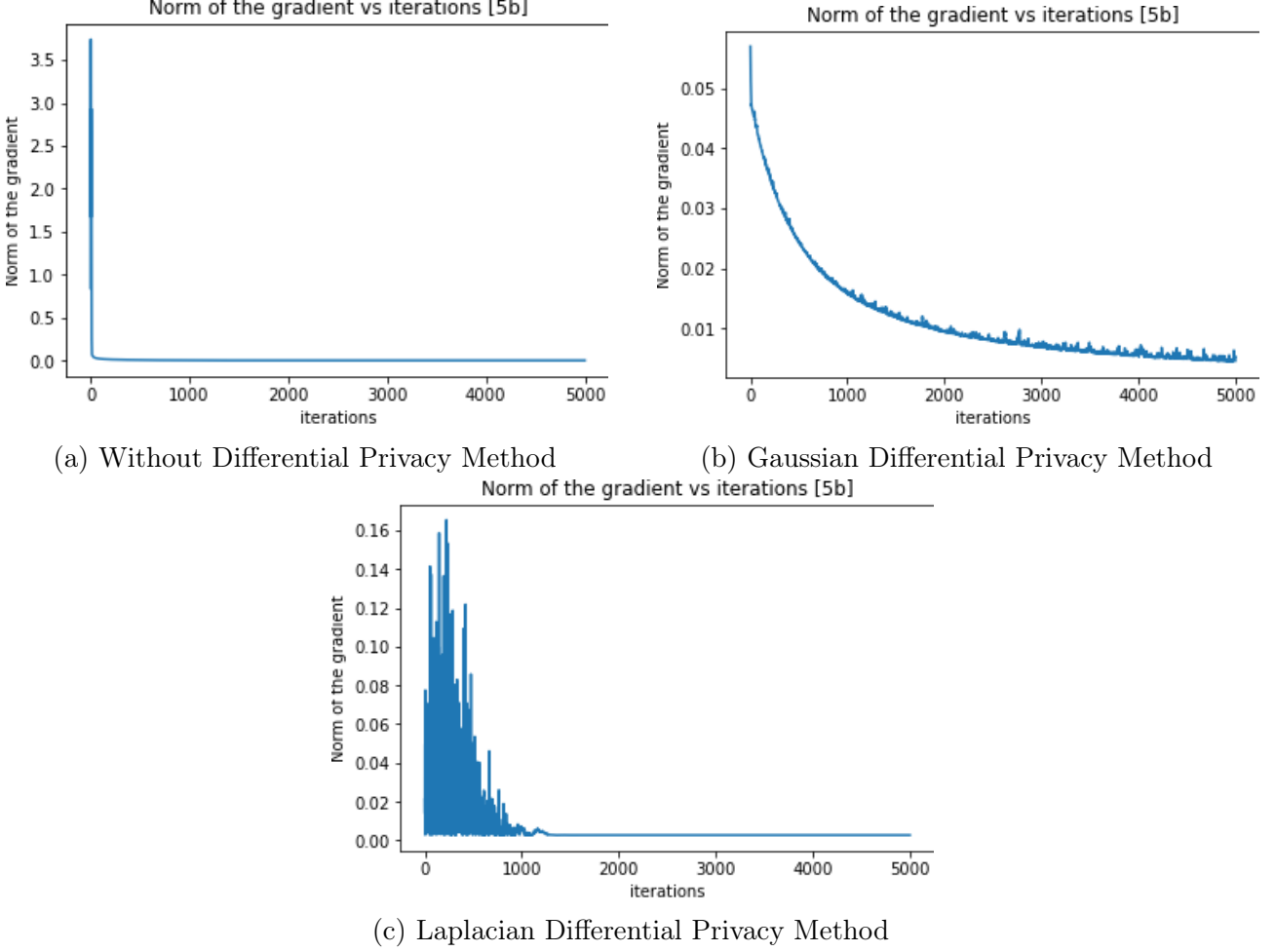


Figure 3: Norm of the gradient vs iterations.

4.3 Precision and recall vs iterations

After that, we evaluate our model on the test dataset by loading the best coefficients. We vary the threshold from 0 to 1 and plot precision and recall versus thresholds. The plots are given in the figure 4. Abrupt behavior is noticed during the evaluation of the model when the laplacian mechanism is applied. The non-differential privacy method has very high precision and recall values in almost all the thresholds. The gaussian method achieves slightly less precision and recall values but they are very comparable to the non-differential method.

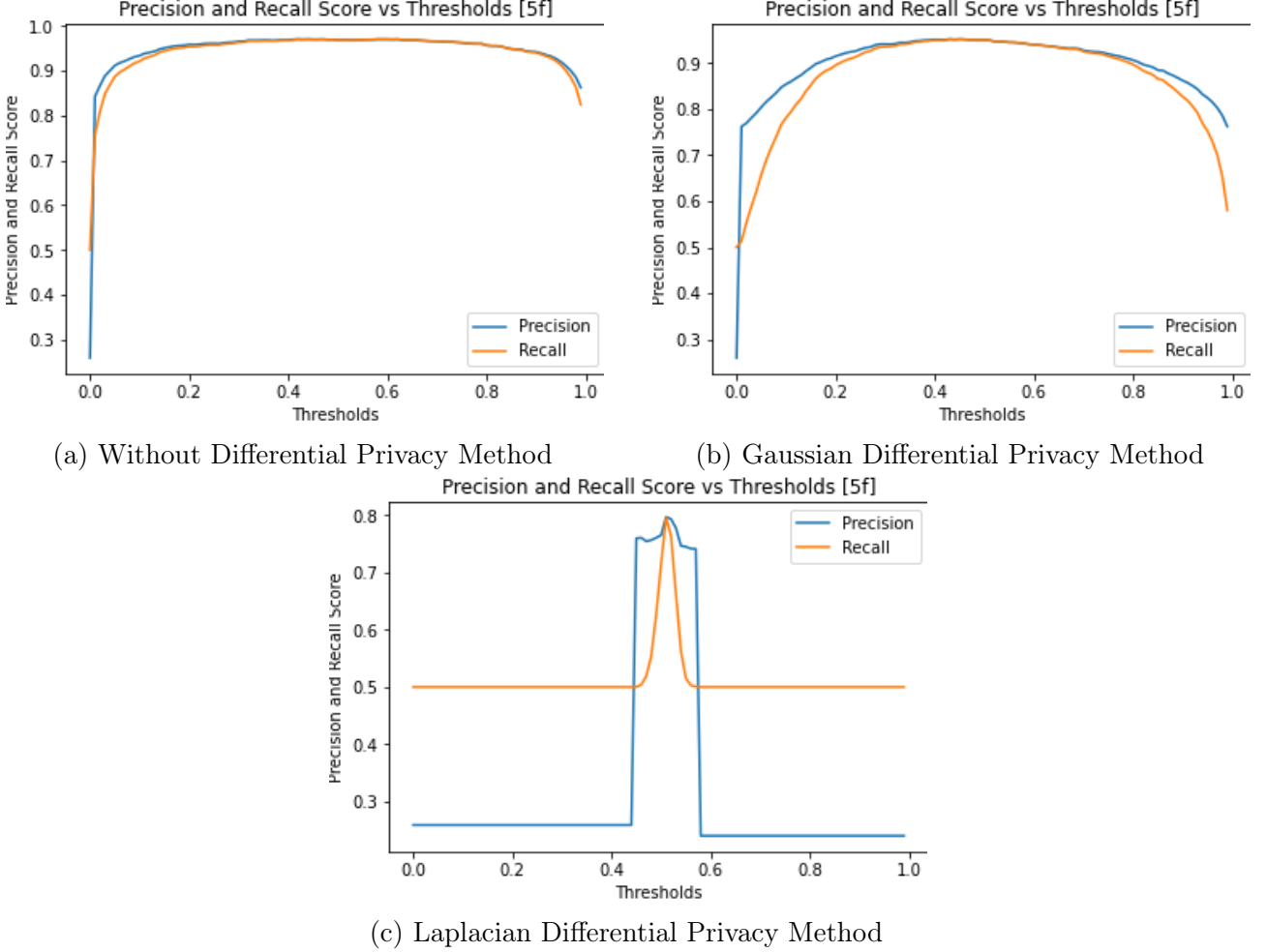


Figure 4: Precision and recall vs iterations.

4.4 True positive rate vs False positive rate

Now, we plot the ROC curve of our model prediction by plotting the true positive rate versus the false positive rate in figure 5. The non-differential method achieved a ROC area of 0.993 while the Gaussian achieved 0.989 and the Laplacian achieved 0.878. It is expected as the Gaussian method usually achieves a better result than the Laplacian method. Also, the non-differential method achieves the highest ROC value due to no privacy mechanism. The privacy mechanisms compromise small performance by giving more security to the data.

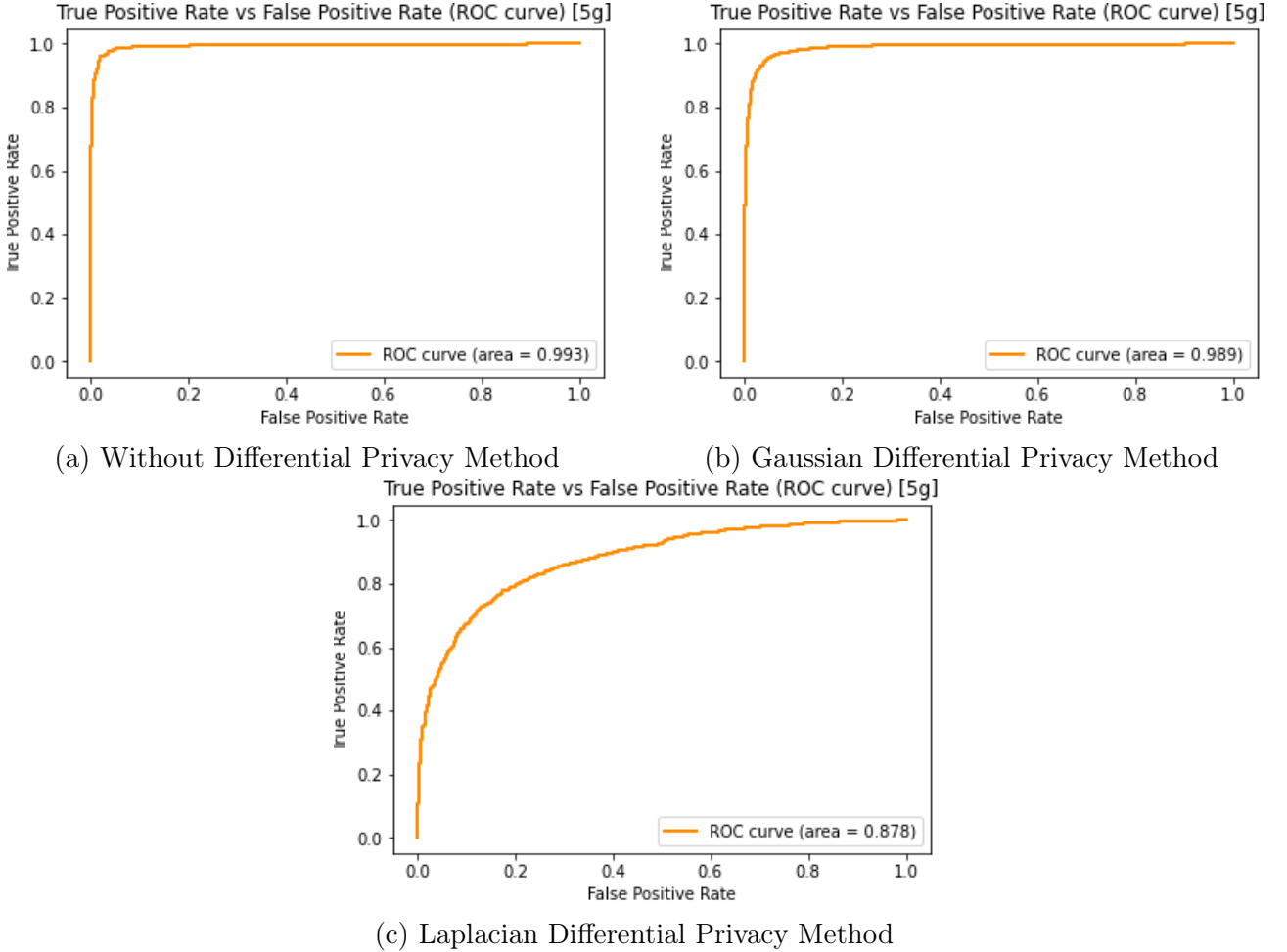


Figure 5: True positive rate vs False positive rate.

4.5 Training and validation loss vs Number of training samples

For demonstrating the effect of the number of training samples on training and validation loss, we gradually decrease the number of training samples keeping the validation dataset intact, and plot training and validation loss vs the number of training samples in figure 6. We can see from the figure that, for the non-differential privacy method as the number of training samples decreases, the validation loss is increasing as well. It is expected as less number of training samples means fewer data available for the model to learn. Inversely, training loss is decreasing at the lower number of training samples because of overfitting. However, for differential privacy methods this relation does not hold. With the increasing number of training samples, both training and validation loss are decreasing. Because, noise is added in each iteration and as a result, the model cannot overfit. As a result, the difference between training and validation loss is very low here.

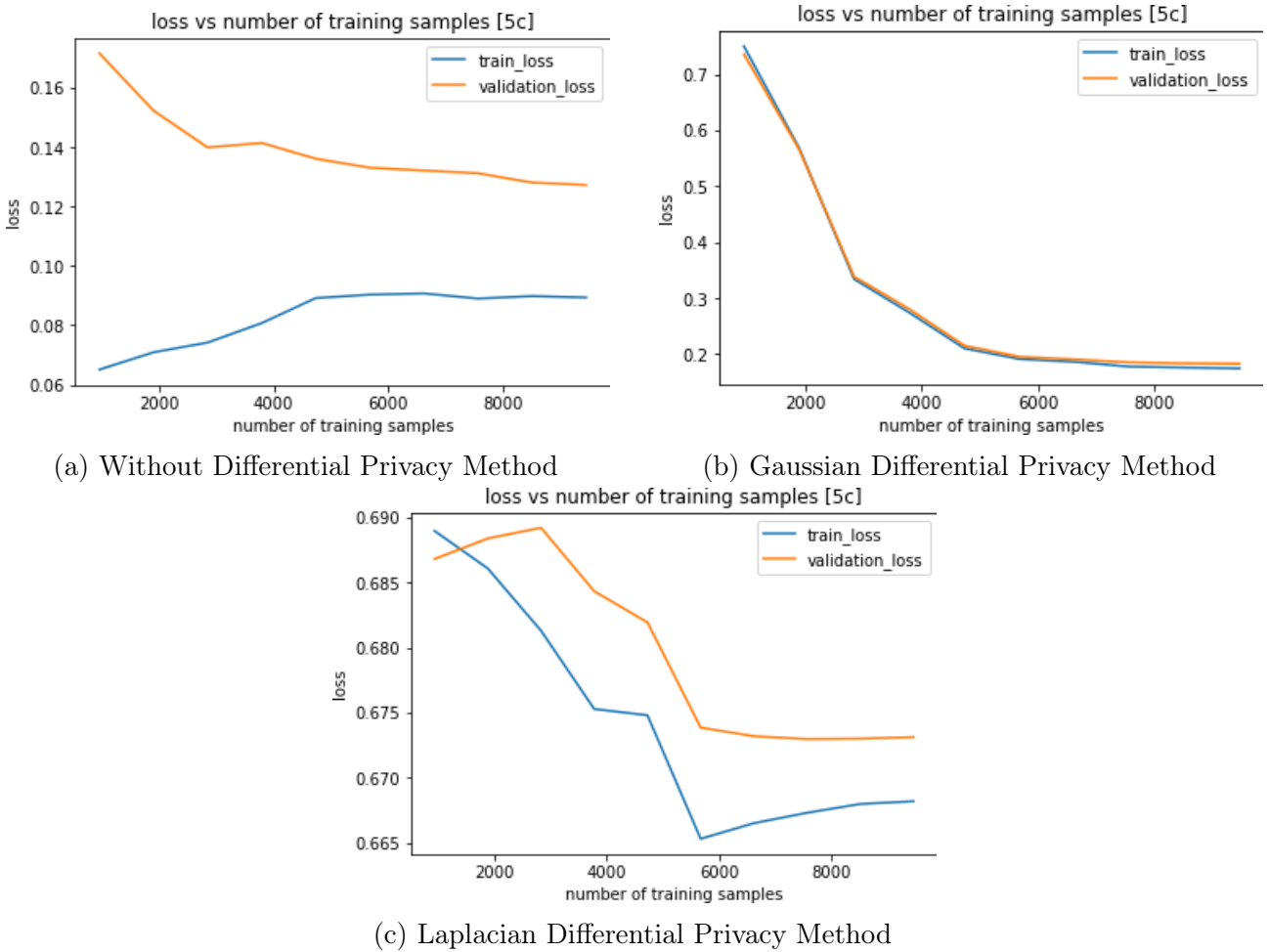


Figure 6: Training and validation loss vs Number of training samples.

The code to calculate these plots are given below:

```
1 def check_num_samples_effect(add_x0=True, iterations=5000, learning_rate
2   =0.5, regularization_coefficient=0.001, method='non_dp',
3   epsilon=0.01, delta=1e-5, data_norm=False, norm_order=2):
4
5   keep_train_percentage = [1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
6   train_losses=[]; val_losses=[]; all_number_training_samples=[];
7   for percentage in keep_train_percentage:
8       print('Current percentage: {}'.format(percentage))
9       X_train, y_train, X_validation, y_validation, X_test, y_test =
10      get_data(data_norm, norm_order)
```

```

9         new_len = int(X_train.shape[0]*percentage)
10        X_train = X_train[0:new_len,...]
11        y_train = y_train[0:new_len]
12        best_coefficients, all_train_loss, all_validation_loss,
all_gradients_norm, _ = logistic_regression(
13                                X_train, y_train,
14                                X_validation, y_validation,
15                                add_x0,
16                                iterations,
17                                learning_rate,
18                                regularization_coefficient,
19                                method,
20                                epsilon,
21                                delta,
22                                )
23
24        train_losses.append(np.min(all_train_loss))
25        val_losses.append(np.min(all_validation_loss))
26        all_number_training_samples.append(new_len)
27        time.sleep(1)
28
29        all_train_losses = np.array(train_losses)
30        all_val_losses = np.array(val_losses)
31        all_number_training_samples = np.array(all_number_training_samples)
32
33        idx = np.argmin(all_val_losses)
34        print('Best val_loss achieved is: {} at num_samples: {}'.format(
all_val_losses[idx], all_number_training_samples[idx]))
35
36        os.makedirs('./logs/', exist_ok=True)
37        # plot of 5d. train and validation loss vs number of training samples
38        plt.plot(all_number_training_samples, all_train_losses, label='
train_loss')
39        plt.plot(all_number_training_samples, all_val_losses, label='
validation_loss')
40        plt.legend()
41        plt.xlabel('number of training samples')
42        plt.ylabel('loss')
43        plt.title('loss vs number of training samples [5c]')
44        plt.savefig("./logs/{}_5c_loss_vs_number_of_training_samples.png".format
(method), bbox_inches='tight', pad_inches=0)
45        plt.show()
46
47        log_dict = {
48            'all_train_loss': all_train_losses,
49            'all_validation_loss': all_val_losses,
50            'all_number_training_samples': all_number_training_samples,
51        }
52
53        np.save('./logs/{}_log_dict_3.npy'.format(method), log_dict)
54    if __name__ == '__main__':
55        # fig c for non_dp, gaussian, laplace
56        check_num_samples_effect(iterations=5000, method='non_dp', data_norm=
False)
57        check_num_samples_effect(iterations=5000, method='gaussian', data_norm=
True, norm_order=2)
58        check_num_samples_effect(iterations=5000, learning_rate=5,
regularization_coefficient=0.001, method='laplace', data_norm=True,
norm_order=1)

```

4.6 Training and validation loss vs Regularization coefficient

We vary the regularization coefficient and calculate the best training and validation loss for each regularization coefficient. We plot these losses vs regularization coefficients in the figure 7. From the figure, we can see that both training and validation loss increases with the regularization coefficient. It is expected as the regularization coefficient add a penalty to the loss function. For very low values of regularization coefficients, the effect on training and validation loss is minimal, almost constant.

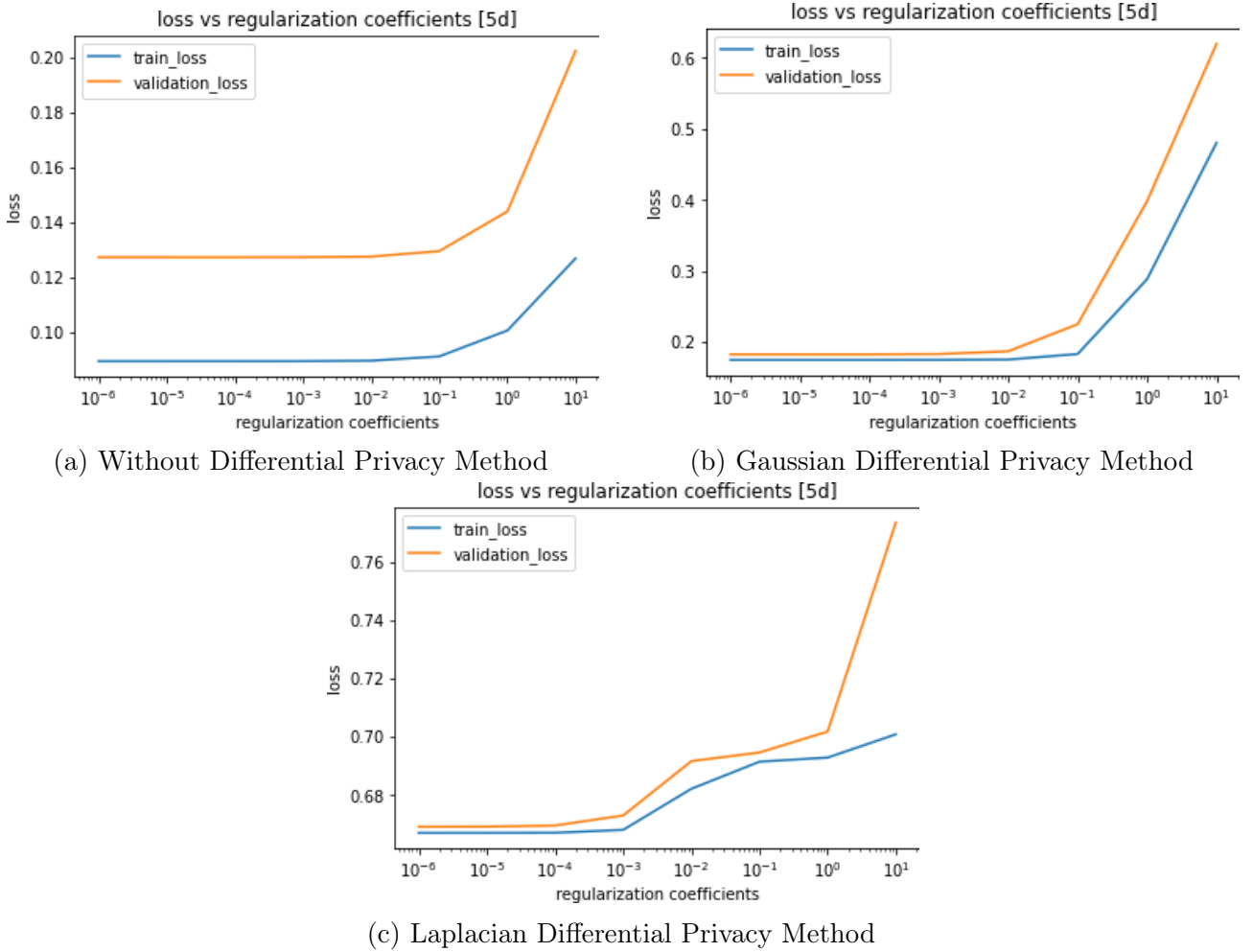


Figure 7: Training and validation loss vs Regularization coefficient.

The code to calculate these plots are given below:

```

1 def check_regularization_effect(add_x0=True, iterations=5000, learning_rate
  =0.5, regularization_coefficient=0.001, method='non_dp',
2     epsilon=0.01, delta=1e-5, data_norm=False, norm_order=2):
3
4     regularization_coefficients=[10, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001,
  0.000001]
5     train_losses=[]; val_losses=[];
6     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data(
  data_norm, norm_order)
7     for regularization_coefficient in regularization_coefficients:
8         print('Current regularization_coefficient: {}'.format(
  regularization_coefficient))
9         best_coefficients, all_train_loss, all_validation_loss,
  all_gradients_norm, _ = logistic_regression(
10             X_train, y_train,
11             X_validation, y_validation,
12             add_x0,
13             iterations,
14             learning_rate,
15             regularization_coefficient,
16             method,
17             epsilon,
18             delta,
19             )
20
21         train_losses.append(np.min(all_train_loss))
22         val_losses.append(np.min(all_validation_loss))
23         time.sleep(1)
24
25     all_train_losses = np.array(train_losses)
26     all_val_losses = np.array(val_losses)
27     regularization_coefficients = np.array(regularization_coefficients)
28     idx = np.argmin(all_val_losses)
29     print('Best val_loss achieved is: {} at regulariztion coefficient: {}'.
  format(all_val_losses[idx], regularization_coefficients[idx]))
30     os.makedirs('./logs/', exist_ok=True)
31     # plot of 5d. train and validation loss vs regularization coefficients
32     plt.plot(regularization_coefficients, all_train_losses, label='
  train_loss')
33     plt.plot(regularization_coefficients, all_val_losses, label='
  validation_loss')
34     plt.xscale('log')
35     plt.legend()
36     plt.xlabel('regularization coefficients')
37     plt.ylabel('loss')
38     plt.title('loss vs regularization coefficients [5d]')
39     plt.savefig("./logs/{}/5d_loss_vs_regularization_coefficients.png".
  format(method), bbox_inches='tight', pad_inches=0)
40     plt.show()
41     log_dict = {
42         'all_train_loss': all_train_losses,
43         'all_validation_loss': all_val_losses,
44         'regularization_coefficients': regularization_coefficients,
45     }
46     np.save('./logs/{}/log_dict_2.npy'.format(method), log_dict)
47 if __name__ == '__main__':
48     # fig d for non_dp, gaussian, laplace
49     check_regularization_effect(iterations=5000, method='non_dp', data_norm=
  False)

```

```

50     check_regularization_effect(iterations=5000, method='gaussian',
data_norm=True, norm_order=2)
51     check_regularization_effect(iterations=5000, learning_rate=5,
regularization_coefficient=0.001, method='laplace', data_norm=True,
norm_order=1)

```

4.7 Required number of iterations vs Step size

Now, we vary the step size and record the number of iterations required for the convergence of the model. We also recorded the best validation loss for each of the step sizes. We plot the required number of iterations vs step size in the figure 8 with blobs representing the best validation losses. The blobs are plotted in the 'jet' color map. Red color blob means highest loss validation value, while royal blue means lowest validation loss value.

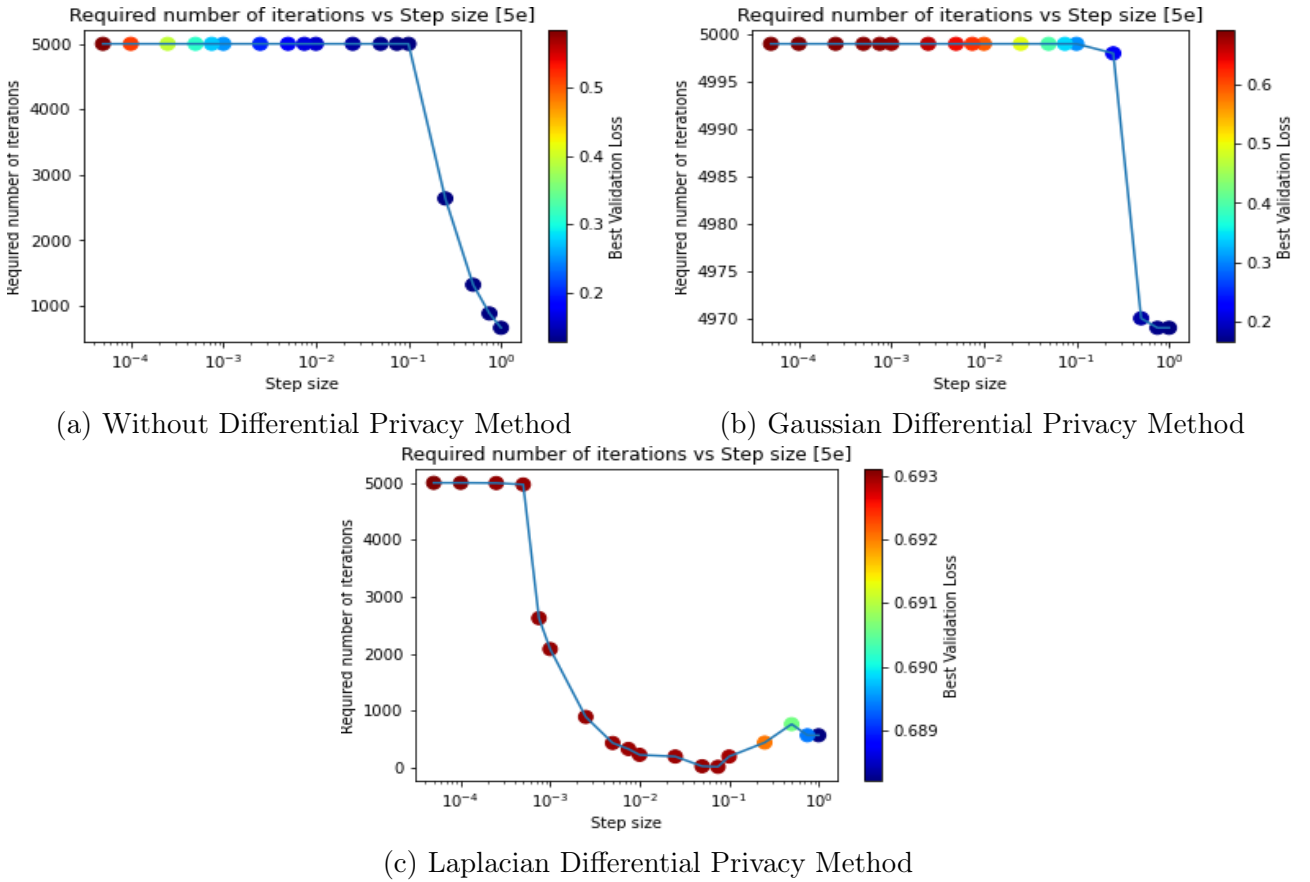


Figure 8: Required number of iterations vs Step size.

From the figure, we can see that the required number of iterations is decreasing with the increase of step size. The best validation loss is decreasing with the increase in the step size. For non-differential privacy and the Gaussian method, the best validation loss is achieved around step size=0.5. The required number of iterations is also less at this step size compared to lower step size values. For the Laplacian method, the best validation loss is achieved at step size=1. It required a bit more iterations than the lowest required number of iterations which is achieved at a lower step size value. The code to calculate these plots are given below:

```

1 def check_step_size_effect(add_x0=True, iterations=5000, learning_rate=0.5,
regularization_coefficient=0.001, method='non_dp', epsilon=0.01, delta=1e
-5, data_norm=False, norm_order=2):
2     learning_rates=[1, 0.75, 0.5, 0.25, 0.1, 0.075, 0.05, 0.025, 0.01,
0.0075, 0.005, 0.0025, 0.001,

```

```

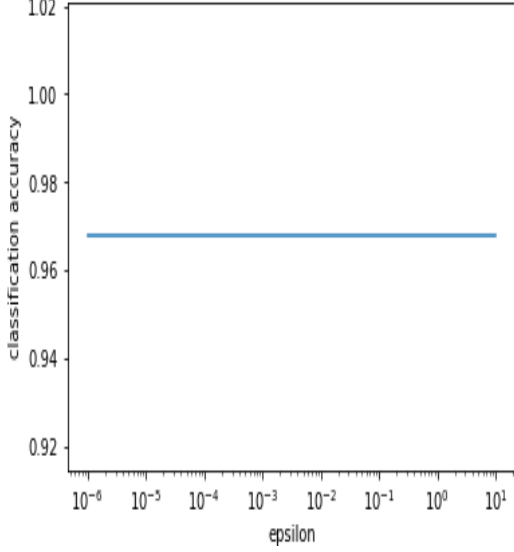
3         0.00075, 0.0005, 0.00025, 0.0001, 0.00005]
4     val_losses=[]; all_required_number_iterations=[];
5     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data(
data_norm, norm_order)
6     for learning_rate in learning_rates:
7         print('Current step size: {}'.format(learning_rate))
8         _, _, all_validation_loss, _, best_iteration = logistic_regression(
9             X_train, y_train,
10            X_validation, y_validation,
11            add_x0,
12            iterations,
13            learning_rate,
14            regularization_coefficient,
15            method,
16            epsilon,
17            delta,
18            )
19         all_required_number_iterations.append(best_iteration)
20         val_losses.append(np.min(all_validation_loss))
21         time.sleep(1)
22     all_required_number_iterations = np.array(all_required_number_iterations
)
23     all_best_loss = np.array(val_losses)
24     step_sizes = np.array(learning_rates)
25     idx = np.argmin(all_best_loss)
26     print('Best val_loss achieved is: {} at step size: {}'.format(
all_best_loss[idx], step_sizes[idx]))
27     os.makedirs('./logs/', exist_ok=True)
28     # plot of 5e. Required number of iterations vs Step size
29     plt.scatter(step_sizes, all_required_number_iterations, s=80, c=
all_best_loss, cmap='jet')
30     plt.plot(step_sizes, all_required_number_iterations)
31     plt.xscale('log')
32     plt.colorbar(label='Best Validation Loss')
33     plt.xlabel('Step size')
34     plt.ylabel('Required number of iterations')
35     plt.title('Required number of iterations vs Step size [5e]')
36     plt.savefig("./logs/{}_5e_Required_number_of_iterations_vs_Step_size.png
".format(method), bbox_inches='tight', pad_inches=0)
37     plt.show()
38     log_dict = {
39         'all_required_number_iterations': all_required_number_iterations,
40         'all_best_loss': all_best_loss,
41         'step_sizes': step_sizes,
42     }
43     np.save('./logs/{}_log_dict_4.npy'.format(method), log_dict)
44 if __name__ == '__main__':
45     # fig e for non_dp, gaussian, laplace
46     check_step_size_effect(iterations=5000, method='non_dp', data_norm=False
)
47     check_step_size_effect(iterations=5000, method='gaussian', data_norm=
True, norm_order=2)
48     check_step_size_effect(iterations=5000, learning_rate=5,
regularization_coefficient=0.001, method='laplace', data_norm=True,
norm_order=1)

```

4.8 Classification accuracy vs epsilon values for fixed delta and sample size.

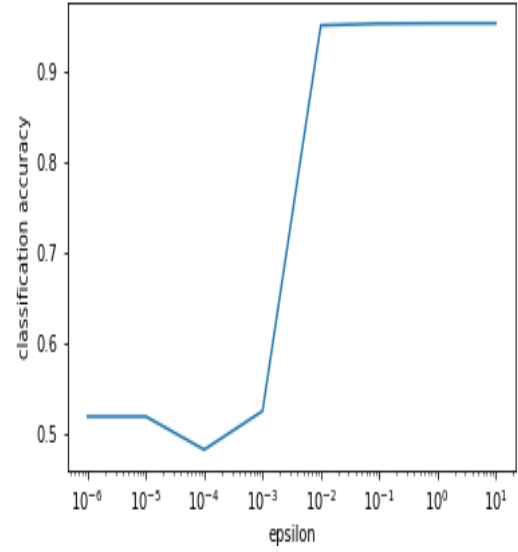
Now, we vary the epsilon values and calculate classification accuracy and ROC value for each of the epsilon values keeping delta and sample size fixed. For the non-differential privacy method, they have no effects as epsilon is only used in differential privacy methods. For differential privacy methods, with the decrease in epsilon value, decrease in classification accuracy and ROC value is observed. A decrease in epsilon value means an increase in privacy. So, the more private the algorithm is, the less accurate the model will be. The plots are shown in the figure 9 10 and code for calculating and plotting them is given afterward.

classification accuracy vs epsilons for sample size=9459 and delta=1e-5 [6]



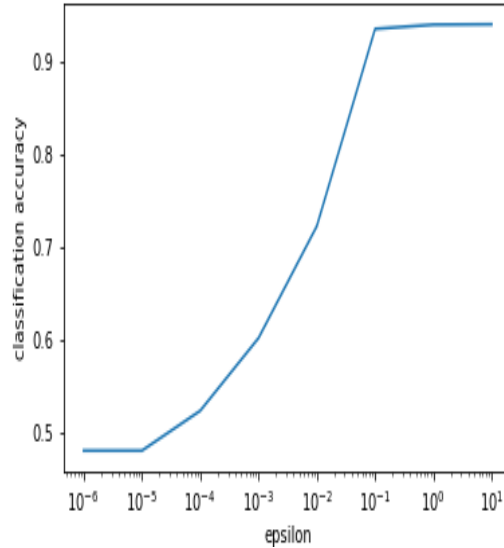
(a) Without Differential Privacy Method

classification accuracy vs epsilons for sample size=9459 and delta=1e-5 [6]



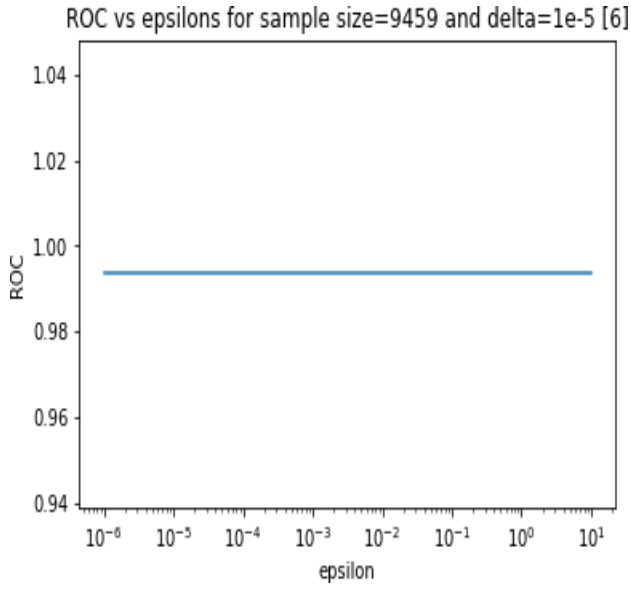
(b) Gaussian Differential Privacy Method

classification accuracy vs epsilons for sample size=9459 and delta=1e-5 [6]

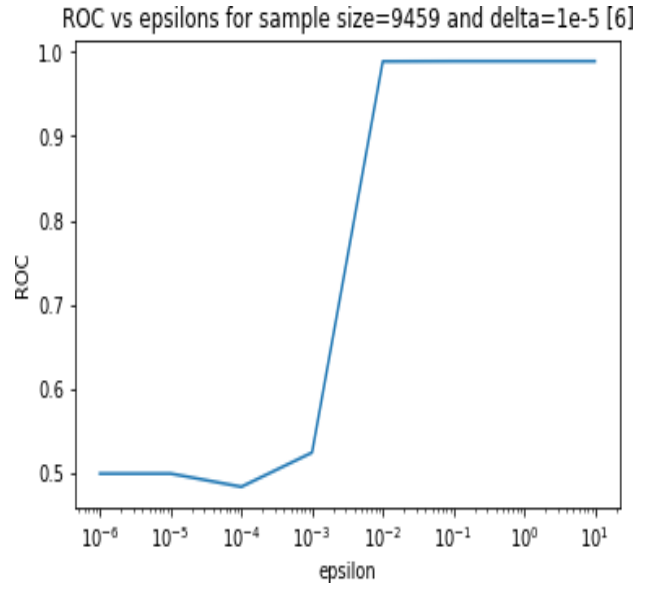


(c) Laplacian Differential Privacy Method

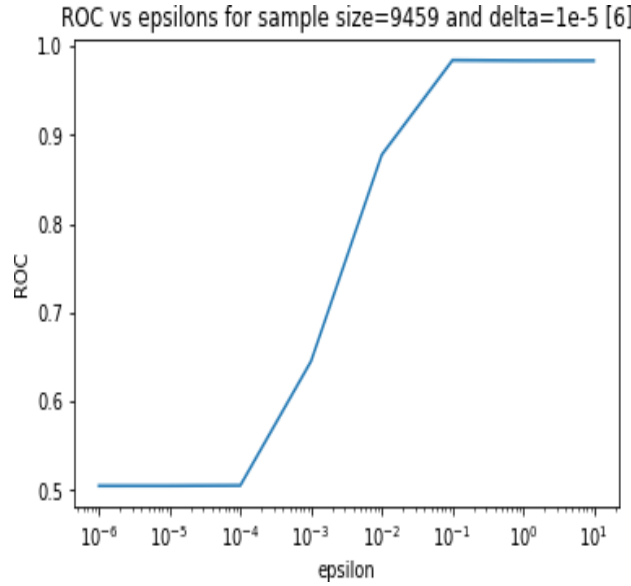
Figure 9: Classification accuracy vs epsilon values for delta=1e-5 and sample size=9459.



(a) Without Differential Privacy Method



(b) Gaussian Differential Privacy Method



(c) Laplacian Differential Privacy Method

Figure 10: ROC vs epsilon values for delta=1e-5 and sample size=9459.

```

1 def check_acc_roc_vs_epsilon_effect(add_x0=True, iterations=5000,
2   learning_rate=0.5, regularization_coefficient=0.001, method='non_dp',
3   epsilon=0.01, delta=1e-5, data_norm=False, norm_order=2):
4
5   all_epsilons = [10, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001, 0.000001]
6   X_train, y_train, X_validation, y_validation, X_test, y_test = get_data(
7     data_norm, norm_order)
8   if add_x0:
9     X_test = np.hstack((np.ones((X_test.shape[0], 1)), X_test))
10
11   all_acc=[];all_roc=[];
12   for epsilon in all_epsilons:
13     best_coefficients, _, _, _, _ = logistic_regression(
14       X_train, y_train,
15       X_validation, y_validation,
16       add_x0,
17       iterations,
18       learning_rate,

```

```

17         regularization_coefficient,
18         method,
19         epsilon,
20         delta,
21     )
22
23     logits = np.dot(X_test, best_coefficients)
24     prediction_probabilites = sigmoid(logits)
25     roc_score = roc_auc_score(y_test, prediction_probabilites, 'macro')
26     y_pred = prediction_probabilites >= 0.5
27     acc = accuracy_score(y_test, y_pred)
28     all_acc.append(acc)
29     all_roc.append(roc_score)
30
31     os.makedirs('./logs/', exist_ok=True)
32     plt.plot(all_epsilons, all_acc)
33     plt.xscale('log')
34     plt.xlabel('epsilon')
35     plt.ylabel('classification accuracy')
36     plt.title('classification accuracy vs epsilons for sample size={} and
37     delta=1e-5 [6]'.format(int(X_train.shape[0])))
38     plt.savefig("./logs/{}_6_classification_accuracy_vs_epsilon.png".format(
39     method), bbox_inches='tight', pad_inches=0)
40     plt.show()
41
42     plt.plot(all_epsilons, all_roc)
43     plt.xscale('log')
44     plt.xlabel('epsilon')
45     plt.ylabel('ROC')
46     plt.title('ROC vs epsilons for sample size={} and delta=1e-5 [6]'.format
47     (int(X_train.shape[0])))
48     plt.savefig("./logs/{}_6_ROC_vs_epsilon.png".format(method), bbox_inches
49     ='tight', pad_inches=0)
50     plt.show()
51
52     log_dict = {
53         'all_acc': np.array(all_acc),
54         'all_epsilons': np.array(all_epsilons),
55         'all_roc': np.array(all_roc)
56     }
57
58     np.save('./logs/{}_log_dict_5.npy'.format(method), log_dict)
59 if __name__ == '__main__':
60     # fig 6 for non_dp, gaussian, laplace
61     check_acc_roc_vs_epsilon_effect(iterations=5000, method='non_dp',
62     data_norm=False)
63     check_acc_roc_vs_epsilon_effect(iterations=5000, method='gaussian',
64     data_norm=True, norm_order=2)
65     check_acc_roc_vs_epsilon_effect(iterations=5000, learning_rate=5,
66     regularization_coefficient=0.001, method='laplace', data_norm=True,
67     norm_order=1)

```

4.9 Classification accuracy vs number of sample sizes for fixed epsilon and delta.

Now, we vary the number of training samples and calculate classification accuracy and ROC value for fixed epsilon and delta value. The plots are given in the figure 11 12. From these figures, it is evident that with an increase in the number of samples, the classification accuracy and ROC value increase for both non-differential and differential privacy methods. Code for calculating and plotting them is given after the figures.

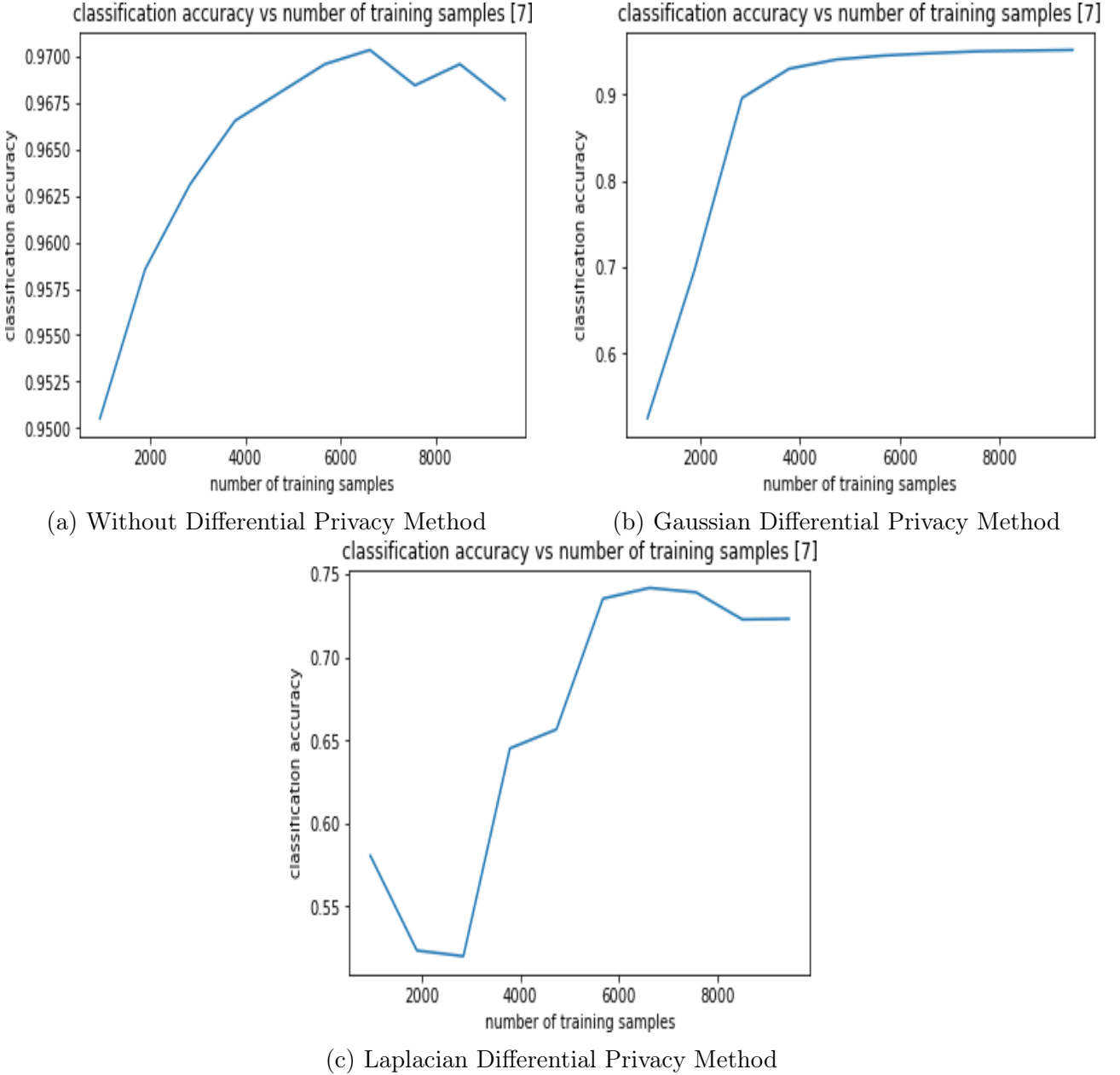
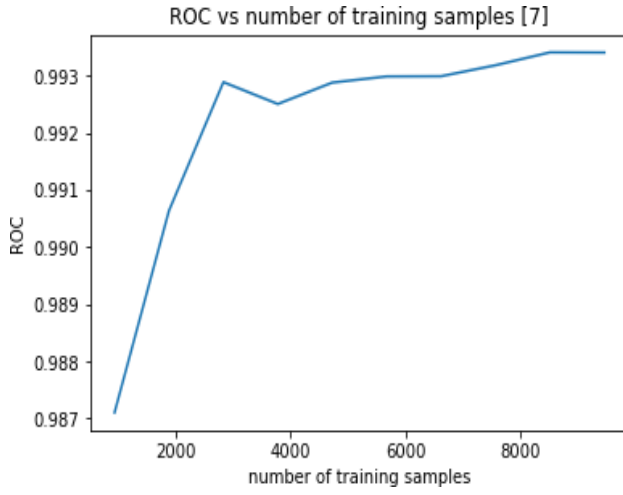
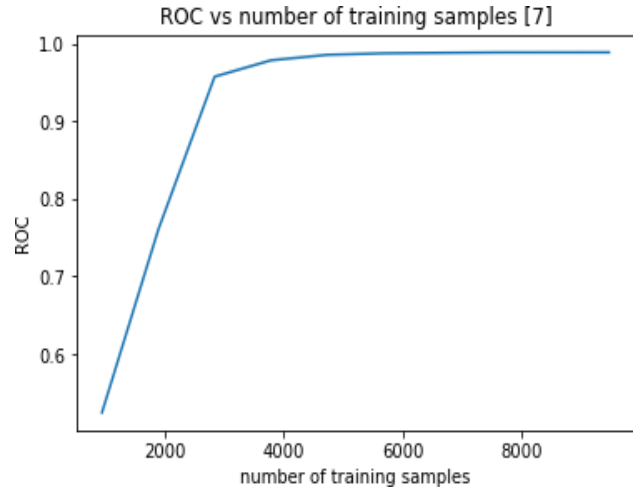


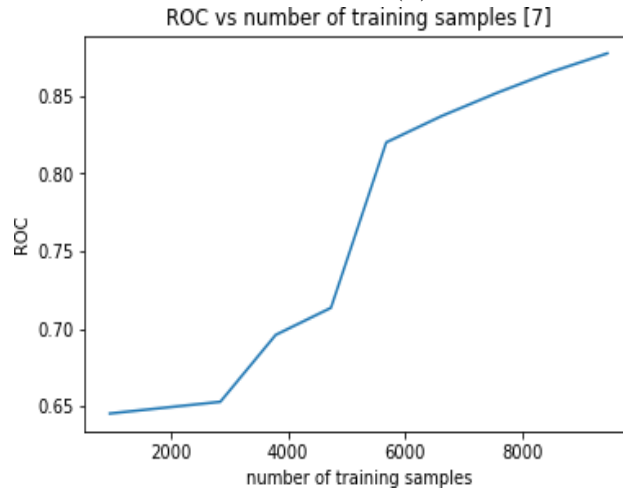
Figure 11: Classification accuracy vs number of sample sizes for epsilon=0.01 and delta=1e-5.



(a) Without Differential Privacy Method



(b) Gaussian Differential Privacy Method



(c) Laplacian Differential Privacy Method

Figure 12: ROC vs number of sample sizes for $\epsilon=0.01$ and $\delta=1e-5$.

```

1 def check_acc_roc_vs_num_samples_effect(add_x0=True, iterations=5000,
2     learning_rate=0.5, regularization_coefficient=0.001,
3     method='non_dp', epsilon=0.01, delta=1e
4     -5, data_norm=False, norm_order=2):
5
6     keep_train_percentage = [1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
7     all_number_training_samples=[];all_acc=[];all_roc=[];
8     for percentage in keep_train_percentage:
9         print('Current percentage: {}'.format(percentage))
10        X_train, y_train, X_validation, y_validation, X_test, y_test =
11        get_data(data_norm, norm_order)
12        new_len = int(X_train.shape[0]*percentage)
13        X_train = X_train[0:new_len,...]
14        y_train = y_train[0:new_len]
15        best_coefficients, _, _, _ = logistic_regression(
16            X_train, y_train,
17            X_validation, y_validation,
18            add_x0,
19            iterations,
20            learning_rate,
21            regularization_coefficient,
22            method,
23            epsilon,
24            delta,

```



```

22         )
23
24     if add_x0:
25         X_test = np.hstack((np.ones((X_test.shape[0], 1)), X_test))
26         logits = np.dot(X_test, best_coefficients)
27         prediction_probabilites = sigmoid(logits)
28         roc_score = roc_auc_score(y_test, prediction_probabilites, 'macro')
29         y_pred = prediction_probabilites >= 0.5
30         acc = accuracy_score(y_test, y_pred)
31         all_acc.append(acc)
32         all_roc.append(roc_score)
33         all_number_training_samples.append(new_len)
34         time.sleep(1)
35
36     all_acc = np.array(all_acc)
37     all_roc = np.array(all_roc)
38     all_number_training_samples = np.array(all_number_training_samples)
39     idx = np.argmax(all_acc)
40     print('Best acc achieved is: {} at num_samples: {}'.format(all_acc[idx],
41         all_number_training_samples[idx]))
42     idx = np.argmax(all_roc)
43     print('Best roc achieved is: {} at num_samples: {}'.format(all_roc[idx],
44         all_number_training_samples[idx]))
45
46     os.makedirs('./logs/', exist_ok=True)
47     # plot of 7 accuracy vs number of training samples
48     plt.plot(all_number_training_samples, all_acc)
49     plt.xlabel('number of training samples')
50     plt.ylabel('classification accuracy')
51     plt.title('classification accuracy vs number of training samples [7]')
52     plt.savefig("./logs/{}_
53     _7_classification_accuracy_vs_number_of_training_samples.png".format(
54     method), bbox_inches='tight', pad_inches=0)
55     plt.show()
56
57     os.makedirs('./logs/', exist_ok=True)
58     # plot of 7 ROC vs number of training samples
59     plt.plot(all_number_training_samples, all_roc)
60     plt.xlabel('number of training samples')
61     plt.ylabel('ROC')
62     plt.title('ROC vs number of training samples [7]')
63     plt.savefig("./logs/{}_
64     _7_roc_vs_number_of_training_samples.png".format(
65     method), bbox_inches='tight', pad_inches=0)
66     plt.show()
67
68     log_dict = {
69         'all_acc': all_acc,
70         'all_roc': all_roc,
71         'all_number_training_samples': all_number_training_samples,
72     }
73     np.save('./logs/{}_log_dict_6.npy'.format(method), log_dict)
74 if __name__ == '__main__':
75     # fig 7 for non_dp, gaussian, laplace
76     check_acc_roc_vs_num_samples_effect(iterations=5000, method='non_dp',
77     data_norm=False)
78     check_acc_roc_vs_num_samples_effect(iterations=5000, method='gaussian',
79     data_norm=True, norm_order=2)
80     check_acc_roc_vs_num_samples_effect(iterations=5000, learning_rate=5,
81     regularization_coefficient=0.001, method='laplace', data_norm=True,
82     norm_order=1)

```