

BANGLADESH UNIVERSITY OF ENGINEERING AND TECHNOLOGY

EEE 6002 Mid-Term 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

Submission Date: 27 February, 2021

1 Assignment Question

Midterm assignment instructions,

- Choose any dataset from the UCI Machine Learning Repository.
- Make sure the dataset satisfies the following criteria:
 - It is a binary classification dataset.
 - The feature dimension is at least 500.
- 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).
- Include the following plots:
 - 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)
- The report must contain a brief description of the dataset and the complete code.

2 Dataset Description

For this assignment, we have selected the following dataset: [Parkinson's Disease Classification Data Set](#). The data used in this study were gathered from 188 patients with PD (107 men and 81 women) with ages ranging from 33 to 87 at the Department of Neurology in CerrahpaÅŸa Faculty of Medicine, Istanbul University. The control group consists of 64 healthy individuals (23 men and 41 women) with ages varying between 41 and 82. During the data collection process, the microphone is set to 44.1 KHz and following the physicians examination, the sustained phonation of the vowel /a/ was collected from each subject with three repetitions.

The total number of data instances in this dataset is 756. Each set of data has a feature dimension of 753. 564 instances has been diagnosed with parkinson disease resulting in class imbalance. The dataset fulfills the criteria set by the assignment. It has at least 500 features per patient and it is a binary classification dataset. Various speech signal processing algorithms including Time Frequency Features, Mel Frequency Cepstral Coefficients (MFCCs), Wavelet Transform based Features, Vocal Fold Features and TWQT features have been applied to the speech recordings of patients to extract these 753 clinically useful information for PD assessment.

First, the dataset is split into 80% and 20% respectively for train and test dataset. This test dataset is used for only evaluation of our trained model. The 80% dataset is further split into 90% and 10% respectively for train and validation dataset. The model is trained on the train dataset and the best coefficients for the model is selected on the result based on validation dataset. Before splitting the dataset, the whole dataset has been normalized into $\{0, 1\}$ range. The data distribution is shown in the table 1.

Table 1: DATA DISTRIBUTION.

Type	Train	Validation	Test
Sample Number	543	61	152
Percentage	72	8	20

The function that I use to process the dataset and required libraries for my implementation are given next:

```
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 sklearn.model_selection import train_test_split
9 from sklearn.metrics import precision_score, recall_score, roc_curve,
    roc_auc_score
10 import pandas as pd
11 from sklearn.preprocessing import MinMaxScaler
12
13 def get_data():
14     # read the dataset
15     df = pd.read_csv('./pd_speech_features.csv')
16     X_features = df.values[1:,-1].astype(np.float)
17     y_labels = df.values[1:,-1].astype(np.int)
```

```

19 # normalize the dataset
20 scaler = MinMaxScaler()
21 scaler = scaler.fit(X_features)
22 X_features = scaler.transform(X_features)
23
24 # train, validation, test split
25 X, X_test, y, y_test = train_test_split(X_features, y_labels, test_size=0.20, random_state=42)
26 X_train, X_validation, y_train, y_validation = train_test_split(X, y, test_size=0.10, random_state=42)
27
28 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. The functions implementing the above algorithm are given below:

```

1 def set_random_state(seed_value):
2     random.seed(seed_value)
3     torch.manual_seed(seed_value)
4     torch.cuda.manual_seed(seed_value)
5     torch.cuda.manual_seed_all(seed_value)
6     os.environ['PYTHONHASHSEED'] = str(seed_value)
7     torch.backends.cudnn.deterministic = True
8     torch.backends.cudnn.benchmark = False
9     np.random.seed(seed_value)
10
11 def sigmoid(x):
12     return 1 / (1 + np.exp(-x))
13
14 def calculate_loss(features, target, coefficients,
15                    regularization_coefficient):
16     logits = np.dot(features, coefficients)

```

```

16     logits_sigmoid = sigmoid(logits)
17     loss = (- np.sum(target*np.log(logits_sigmoid)+ (1-target)*np.log(1-
18         logits_sigmoid)) \
19         + 0.5*regularization_coefficient*np.dot(coefficients,
20         coefficients))/len(target)
21     return loss
22
23 def train_one_iteration(features, target, coefficients, learning_rate,
24     regularization_coefficient):
25     f_x = np.dot(features, coefficients)
26     p_x = sigmoid(f_x)
27     gradient = (np.dot(features.T, (p_x - target)) +
28         regularization_coefficient * coefficients)/features.shape[0]
29     coefficients -= learning_rate * gradient
30     return coefficients, gradient
31
32 def logistic_regression(features_train, target_train, features_validation,
33     target_validation, add_x0,
34     iterations, learning_rate,
35     regularization_coefficient):
36     set_random_state(42)
37
38     if add_x0:
39         x0 = np.ones((features_train.shape[0], 1))
40         features_train = np.hstack((x0, features_train))
41         x0 = np.ones((features_validation.shape[0], 1))
42         features_validation = np.hstack((x0, features_validation))
43         coefficients = np.zeros(features_train.shape[1])
44
45     best_coefficients = np.zeros(features_train.shape[1])
46     best_loss = 9999
47     count = 0
48     patience = 100
49     best_iteration = 0
50
51     all_train_loss=[]
52     all_validation_loss=[]
53     all_gradients=[]
54     progress_bar = tqdm(total=iterations, unit='iterations')
55     for step in range(iterations):
56         coefficients, gradient = train_one_iteration(features_train,
57             target_train, coefficients, learning_rate, regularization_coefficient)
58
59         train_loss = calculate_loss(features_train, target_train,
60             coefficients, regularization_coefficient)
61         validation_loss = calculate_loss(features_validation,
62             target_validation, coefficients, regularization_coefficient)
63
64         if validation_loss < best_loss:
65             best_loss = validation_loss
66             best_coefficients = coefficients.copy()
67             count = 0
68             best_iteration = step
69         else:
70             count += 1
71
72         progress_bar.set_postfix({'Train loss':train_loss, ' Validation_loss '
73             :validation_loss})
74         progress_bar.update(1)

```

```

66     all_train_loss.append(train_loss)
67     all_validation_loss.append(validation_loss)
68     all_gradients.append(gradient)
69
70     if count == patience:
71         learning_rate = learning_rate * 0.5
72         count = 0
73
74     progress_bar.close()
75
76     all_gradients = np.stack(all_gradients)
77     all_gradients_norm = np.linalg.norm(all_gradients, axis=-1)
78
79     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 curves. For this, we selected step size=0.04, regularization coefficient=0.001 and ran the training for 4000 iterations. 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):
3
4     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data
5     ()
6     best_coefficients, all_train_loss, all_validation_loss,
7     all_gradients_norm, _ = logistic_regression(
8         X_train, y_train,
9         X_validation, y_validation,
10         add_x0,
11         iterations,
12         learning_rate,
13         regularization_coefficient,
14     )
15
16     os.makedirs('./logs/', exist_ok=True)
17     # plot of 5a. train and validation loss vs iterations
18     iterations = np.arange(0, iterations)
19     plt.plot(iterations, all_train_loss, label='train_loss')
20     plt.plot(iterations, all_validation_loss, label='validation_loss')
21     plt.legend()
22     plt.xlabel('iterations')
23     plt.ylabel('loss')
24     plt.title('loss vs iterations [5a]')
25     plt.savefig("./logs/5a_loss_vs_iterations.png", bbox_inches='tight',
26         pad_inches=0)
27     plt.show()
28
29     # plot of 5b. Norm of the gradient vs iterations
30     plt.plot(iterations, all_gradients_norm)
31     plt.xlabel('iterations')
32     plt.ylabel('Norm of the gradient')
33     plt.title('Norm of the gradient vs iterations [5b]')
34     plt.savefig("./logs/5b_Norm_of_the_gradient_vs_iterations.png",
35         bbox_inches='tight', pad_inches=0)
36     plt.show()

```

```

32
33     if add_x0:
34         X_test = np.hstack((np.ones((X_test.shape[0], 1)), X_test))
35     logits = np.dot(X_test, best_coefficients)
36     prediction_probabilites = sigmoid(logits)
37
38     all_precisions=[]
39     all_recalls=[]
40     all_thresholds=np.arange(0, 1, 0.01)
41     for threshold in all_thresholds:
42         y_true = y_test.copy()
43         y_pred = prediction_probabilites >= threshold
44         all_precisions.append(precision_score(y_true, y_pred, average='macro
'))
45         all_recalls.append(recall_score(y_true, y_pred, average='macro'))
46
47     # plot of 5f. Precision and Recall vs threshold
48     plt.plot(all_thresholds, all_precisions, label='Precision')
49     plt.plot(all_thresholds, all_recalls, label='Recall')
50     plt.legend()
51     plt.xlabel('Thresholds')
52     plt.ylabel('Precision and Recall Score')
53     plt.title('Precision and Recall Score vs Thresholds [5f]')
54     plt.savefig("./logs/5f_Precision_and_Recall_Score_vs_Thresholds.png",
bbox_inches='tight', pad_inches=0)
55     plt.show()
56
57     # plot of 5g. True Positive Rate vs False Positive Rate (ROC curve)
58     fpr, tpr, _ = roc_curve(y_true, prediction_probabilites)
59     roc_score = roc_auc_score(y_true, prediction_probabilites, 'macro')
60     plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area =
 {:.3f})'.format(roc_score))
61     plt.legend(loc="lower right")
62     plt.xlabel('False Positive Rate')
63     plt.ylabel('True Positive Rate')
64     plt.title('True Positive Rate vs False Positive Rate (ROC curve) [5g]')
65     plt.savefig("./logs/5g_True_Positive_Rate_vs_False_Positive_Rate.png",
bbox_inches='tight', pad_inches=0)
66     plt.show()
67
68     log_dict = {
69         'all_train_loss': all_train_loss,
70         'all_validation_loss': all_validation_loss,
71         'all_gradients_norm': all_gradients_norm,
72         'iterations': iterations,
73         'all_thresholds': all_thresholds,
74         'all_precisions': all_precisions,
75         'all_recalls': all_recalls,
76         'fpr': fpr,
77         'tpr': tpr,
78         'roc_score': roc_score,
79     }
80
81     np.save('./logs/log_dict_1.npy', log_dict)
82 if __name__ == '__main__':
83     # fig a,b,f,g
84     train(iterations=4000, learning_rate=0.04, regularization_coefficient
=0.001)

```

4.1 Training and Validation loss vs iterations

From the figure 1, we can see that validation loss converges between 3500 and 4000 iterations. The difference between training and validation loss indicates the overfitting of the model.

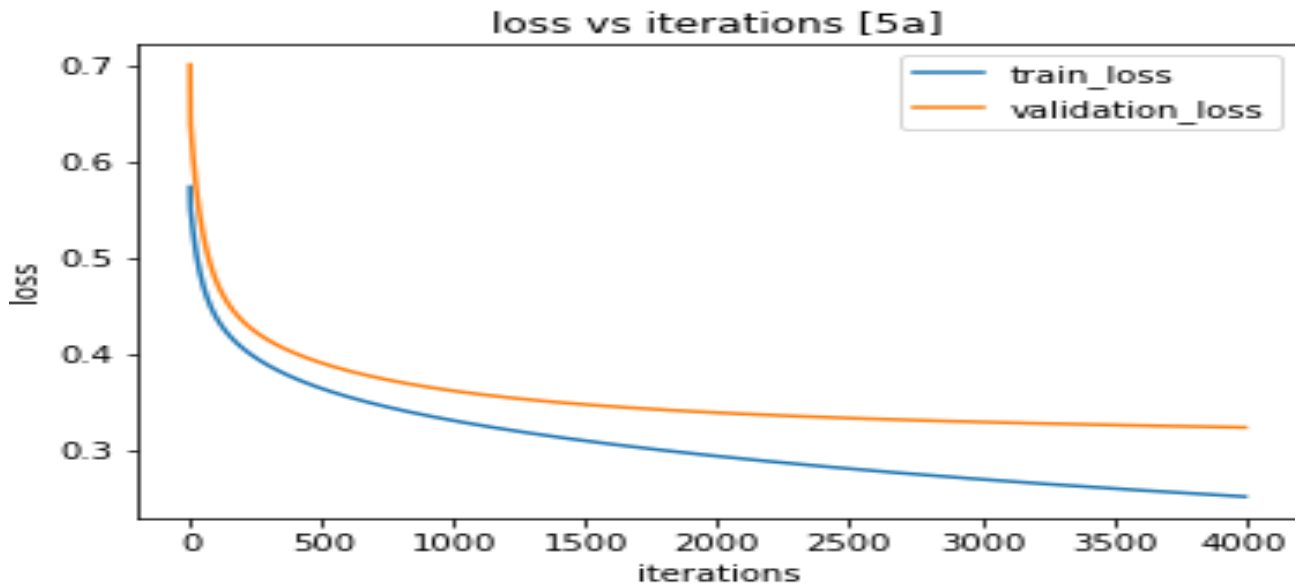


Figure 1: Training and Validation loss vs iterations.

4.2 Norm of the gradient vs iterations

Now, we plot the Norm of the gradient against iterations in the figure 2. 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.

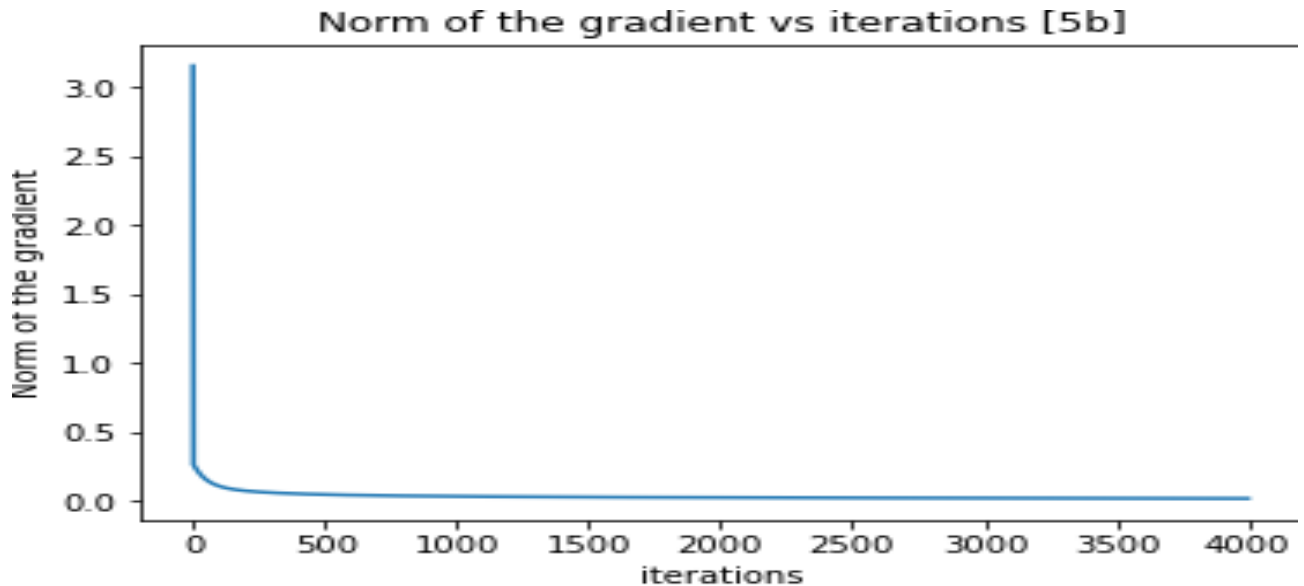


Figure 2: Norm of the gradient vs iterations.

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 3. We can see from the figure that high precision is achieved at lower threshold while high recall is achieved at higher threshold.

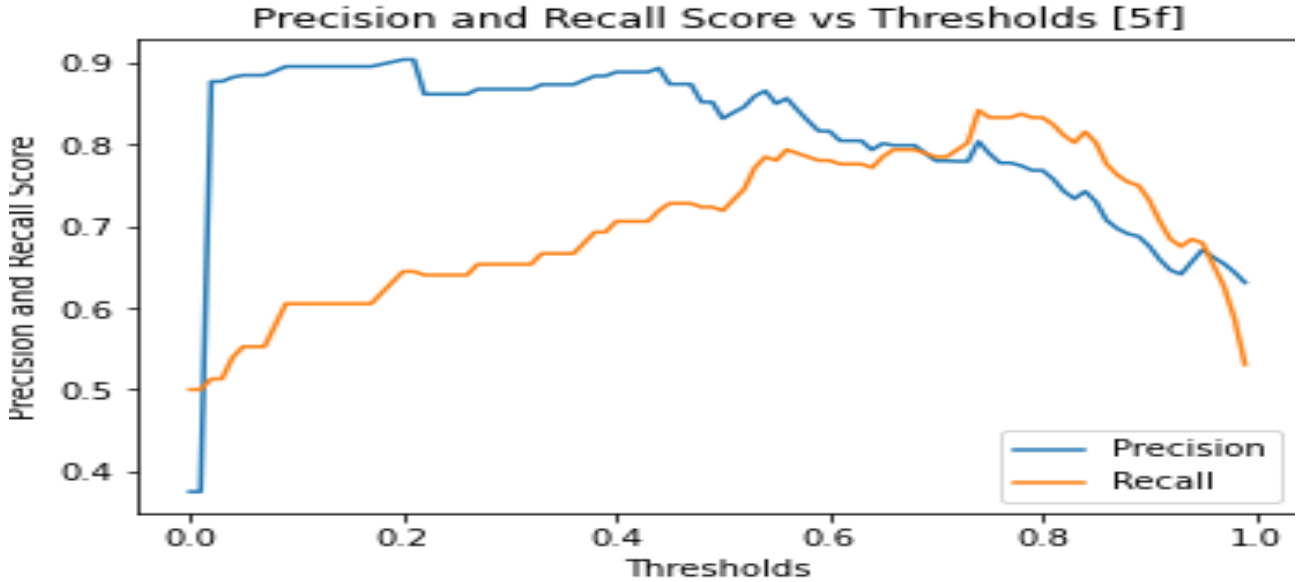


Figure 3: Precision and recall vs iterations.

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 4. Our model has achieved ROC area of 0.891 which is very close to ideal area of 1.

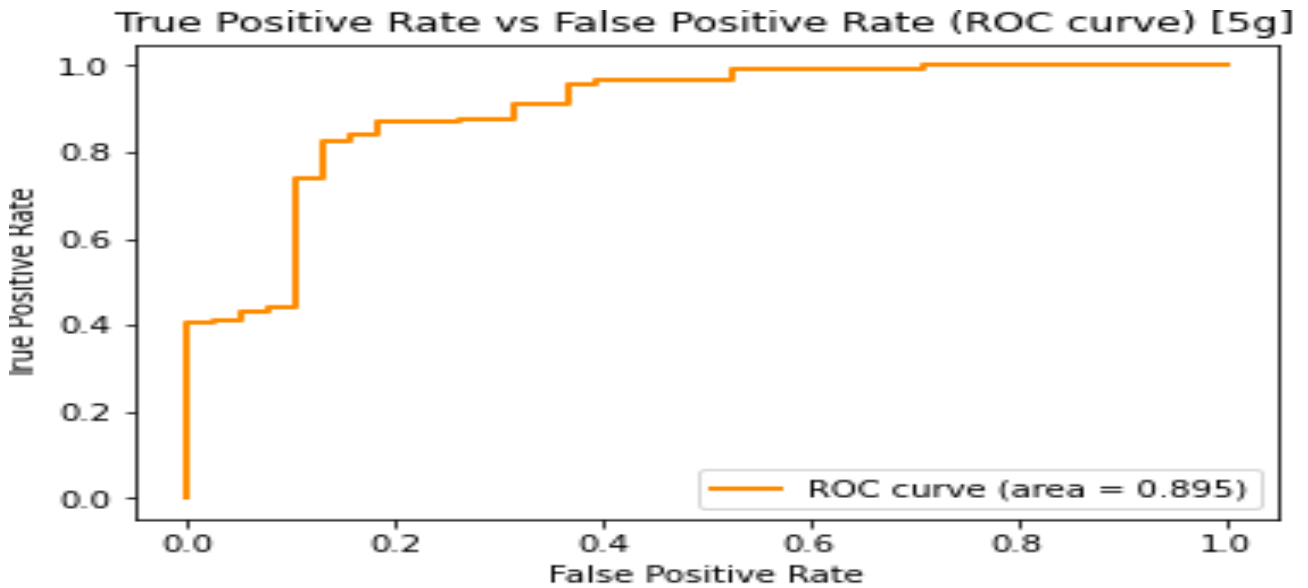


Figure 4: True positive rate vs False positive rate.

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 5. 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.

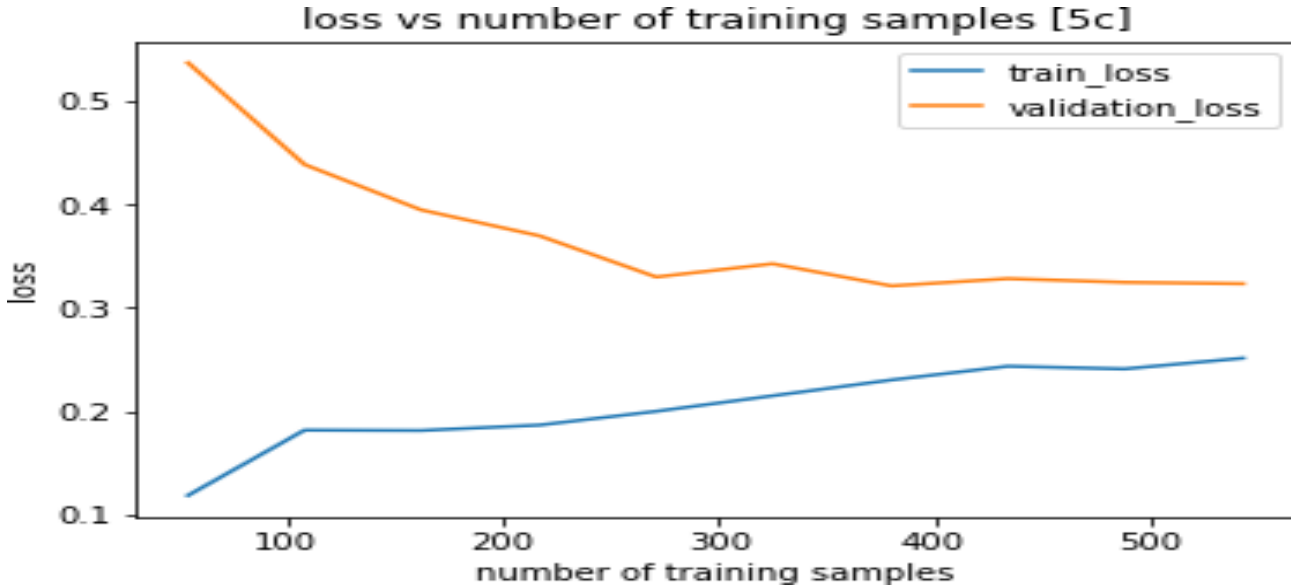


Figure 5: Training and validation loss vs Number of training samples.

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):
3     keep_train_percentage = [1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
4     train_losses=[]; val_losses=[]; all_number_training_samples=[];
5     for percentage in keep_train_percentage:
6         print('Current percentage: {}'.format(percentage))
7         X_train, y_train, X_validation, y_validation, X_test, y_test =
8         get_data()
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,
13        all_gradients_norm, _ = logistic_regression(
14            X_train, y_train,
15            X_validation, y_validation,
16            add_x0,
17            iterations,
18            learning_rate,
19            regularization_coefficient,
20        )
21        train_losses.append(np.min(all_train_loss))
22        val_losses.append(np.min(all_validation_loss))
23        all_number_training_samples.append(new_len)
24        time.sleep(1)
25    all_train_losses = np.array(train_losses)
```

```

25     all_val_losses = np.array(val_losses)
26     all_number_training_samples = np.array(all_number_training_samples)
27     idx = np.argmin(all_val_losses)
28     print('Best val_loss achieved is: {} at num_samples: {}'.format(
all_val_losses[idx], all_number_training_samples[idx]))
29     os.makedirs('./logs/', exist_ok=True)
30     # plot of 5d. train and validation loss vs number of training samples
31     plt.plot(all_number_training_samples, all_train_losses, label='
train_loss')
32     plt.plot(all_number_training_samples, all_val_losses, label='
validation_loss')
33     plt.legend()
34     plt.xlabel('number of training samples')
35     plt.ylabel('loss')
36     plt.title('loss vs number of training samples [5c]')
37     plt.savefig("./logs/5c_loss_vs_number_of_training_samples.png",
bbox_inches='tight', pad_inches=0)
38     plt.show()
39     log_dict = {
40         'all_train_loss': all_train_losses,
41         'all_validation_loss': all_val_losses,
42         'all_number_training_samples': all_number_training_samples,
43     }
44     np.save('./logs/log_dict_3.npy', log_dict)
45 if __name__ == '__main__':
46     # fig c
47     check_num_samples_effect(iterations=4000, learning_rate=0.04,
regularization_coefficient=0.001)

```

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 6. 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.

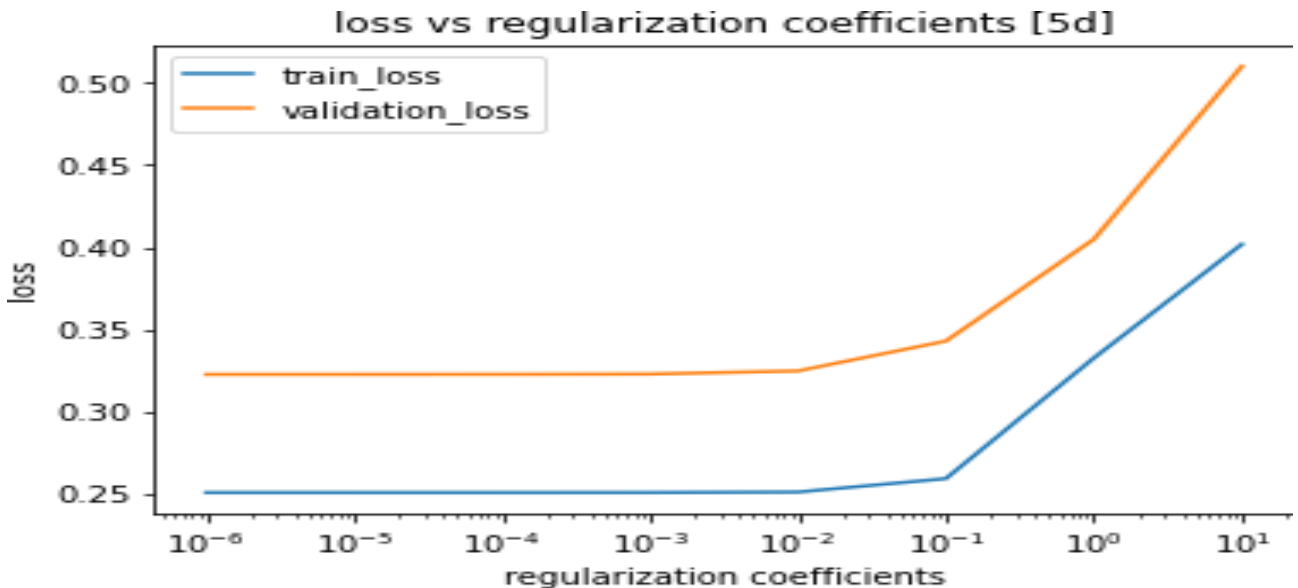


Figure 6: Training and validation loss vs Regularization coefficient.

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):
2
3     regularization_coefficients=[10, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001,
  0.000001]
4     train_losses=[]; val_losses=[];
5     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data
  ()
6     for regularization_coefficient in regularization_coefficients:
7         print('Current regularization_coefficient: {}'.format(
  regularization_coefficient))
8         best_coefficients, all_train_loss, all_validation_loss,
  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            )
16
17         train_losses.append(np.min(all_train_loss))
18         val_losses.append(np.min(all_validation_loss))
19         time.sleep(1)
20
21     all_train_losses = np.array(train_losses)
22     all_val_losses = np.array(val_losses)
23     regularization_coefficients = np.array(regularization_coefficients)
24
25     idx = np.argmin(all_val_losses)
26     print('Best val_loss achieved is: {} at regulariztion coefficient: {}'.
  format(all_val_losses[idx], regularization_coefficients[idx]))
27
28     os.makedirs('./logs/', exist_ok=True)
29     # plot of 5d. train and validation loss vs regularization coefficients
30     plt.plot(regularization_coefficients, all_train_losses, label='
  train_loss')
31     plt.plot(regularization_coefficients, all_val_losses, label='
  validation_loss')
32     plt.xscale('log')
33     plt.legend()
34     plt.xlabel('regularization coefficients')
35     plt.ylabel('loss')
36     plt.title('loss vs regularization coefficients [5d]')
37     plt.savefig("./logs/5d_loss_vs_regularization_coefficients.png",
  bbox_inches='tight', pad_inches=0)
38     plt.show()
39
40     log_dict = {
41         'all_train_loss': all_train_losses,
42         'all_validation_loss': all_val_losses,
43         'regularization_coefficients': regularization_coefficients,
44     }
45
46     np.save('./logs/log_dict_2.npy', log_dict)
47 if __name__ == '__main__':
48     # fig d
49     check_regularization_effect(iterations=4000, learning_rate=0.04,
  regularization_coefficient=0.001)
```

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 7 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.

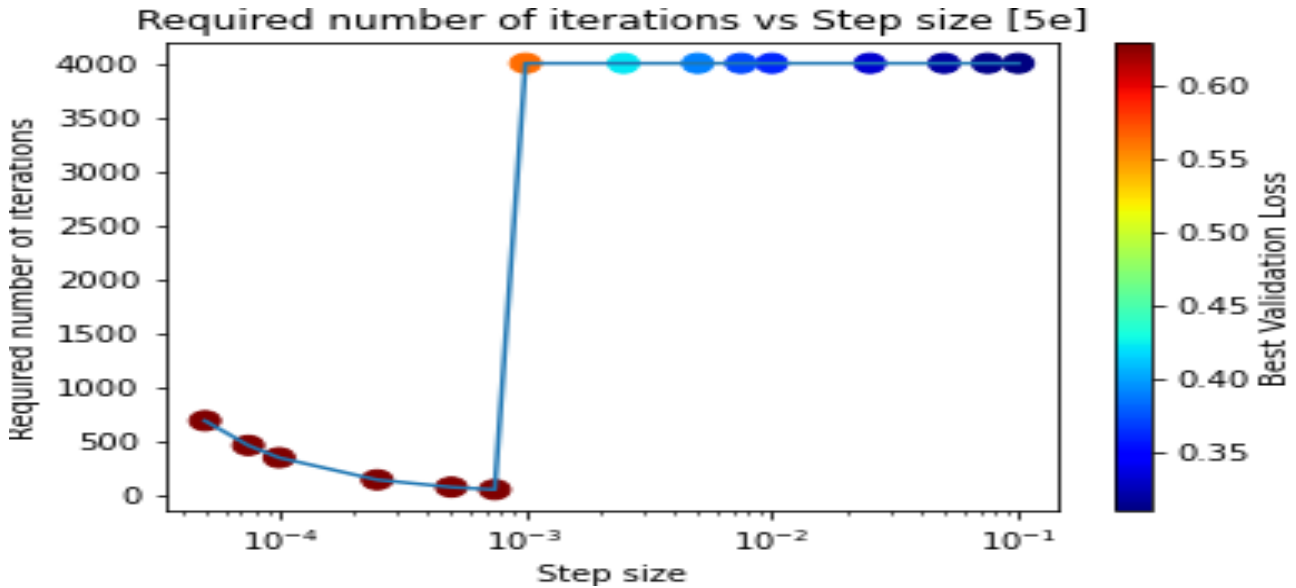


Figure 7: Required number of iterations vs Step size.

```
1 def check_step_size_effect(add_x0=True, iterations=5000, learning_rate=0.5,
2   regularization_coefficient=0.001):
3
4     learning_rates=[0.1, 0.075, 0.05, 0.025, 0.01, 0.0075, 0.005, 0.0025,
5     0.001, 0.00075, 0.0005, 0.00025, 0.0001,
6     0.000075, 0.00005]
7     val_losses=[]; all_required_number_iterations=[];
8     X_train, y_train, X_validation, y_validation, X_test, y_test = get_data
9     ()
10    for learning_rate in learning_rates:
11        print('Current step size: {}'.format(learning_rate))
12        _, _, all_validation_loss, _, best_iteration = logistic_regression(
13            X_train, y_train,
14            X_validation, y_validation,
15            add_x0,
16            iterations,
17            learning_rate,
18            regularization_coefficient,
19            )
20
21        all_required_number_iterations.append(best_iteration)
22        val_losses.append(np.min(all_validation_loss))
23        time.sleep(1)
24
25    all_required_number_iterations = np.array(all_required_number_iterations
26    )
27    all_best_loss = np.array(val_losses)
28    step_sizes = np.array(learning_rates)
29
30    idx = np.argmin(all_best_loss)
```

```

27     print('Best val_loss achieved is: {} at step size: {}'.format(
all_best_loss[idx], step_sizes[idx]))
28
29     os.makedirs('./logs/', exist_ok=True)
30     # plot of 5e. Required number of iterations vs Step size
31     plt.scatter(step_sizes, all_required_number_iterations, s=80, c=
all_best_loss, cmap='jet')
32     plt.plot(step_sizes, all_required_number_iterations)
33     plt.xscale('log')
34     plt.colorbar(label='Best Validation Loss')
35     plt.xlabel('Step size')
36     plt.ylabel('Required number of iterations')
37     plt.title('Required number of iterations vs Step size [5e]')
38     plt.savefig("./logs/5e_Required_number_of_iterations_vs_Step_size.png",
bbox_inches='tight', pad_inches=0)
39     plt.show()
40
41     log_dict = {
42         'all_required_number_iterations': all_required_number_iterations,
43         'all_best_loss': all_best_loss,
44         'step_sizes': step_sizes,
45     }
46
47     np.save('./logs/log_dict_4.npy', log_dict)
48 if __name__ == '__main__':
49     # fig e
50     check_step_size_effect(iterations=4000, learning_rate=0.04,
regularization_coefficient=0.001)

```