# Linear Regression (HW1)

Student:      張皓雲

Student Id:   310551031

# Content

# ● Part. 1, Coding (60%):

**1.** **(15%) Plot the learning curve of the training, you should find that loss decreases after a few iterations (x-axis=iteration, y-axis=loss, Matplotlib or other plot tools is available to use)**

I used the gradient descent method, mini-batch gradient descent method and stohastic gradient descent method to finish this assignment. The loss curves of the three methods are as follows. Obviously, because only one data point is seen in each iteration of stochastic gradient descent, the loss curve of stochastic gradient descent is more oscillating than that of minibatch gradient descent. The loss curve figure, data point with fitting line figure are packing into the compress file.

## A. Gradient Descent Loss Curve
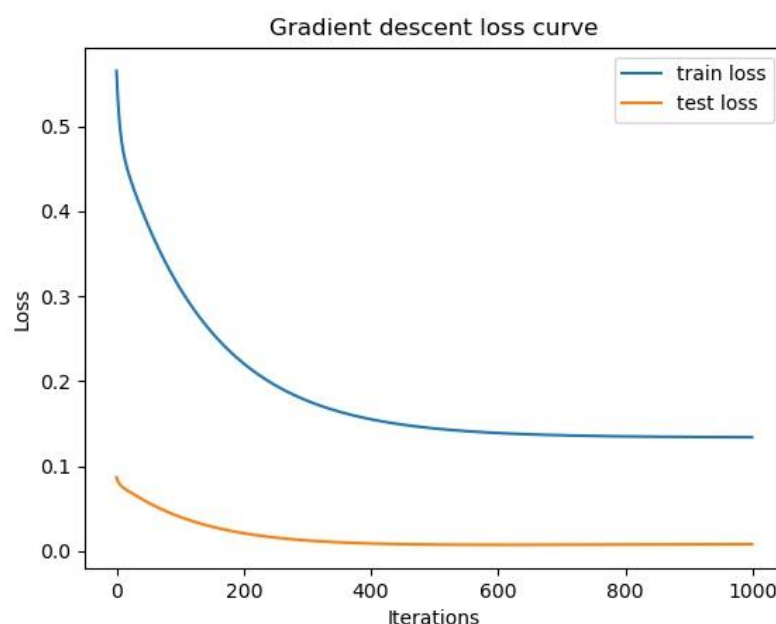


Figure 1. Gradient Descent Loss Curve

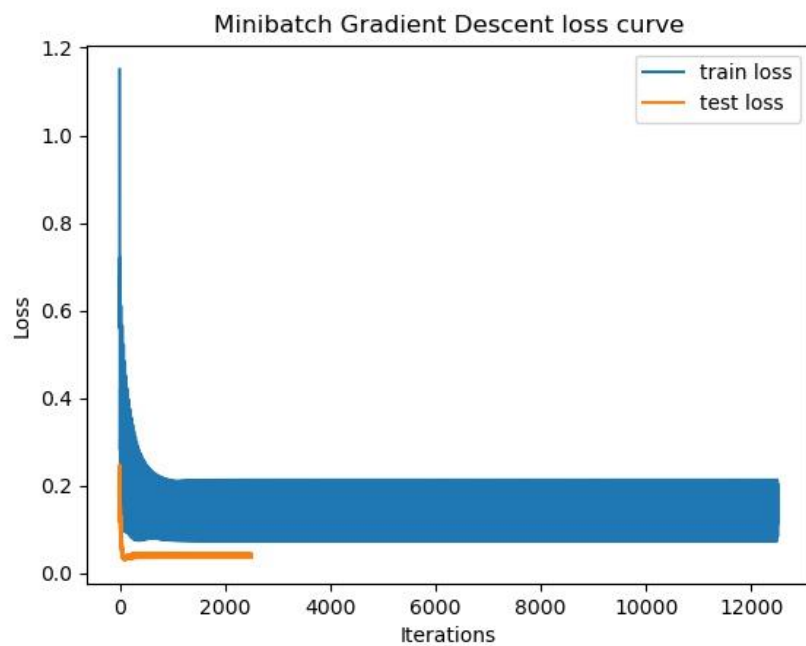## B. Minibatch Gradient Descent Loss Curve



Figure 2. Minibatch Gradient Descent Loss Curve
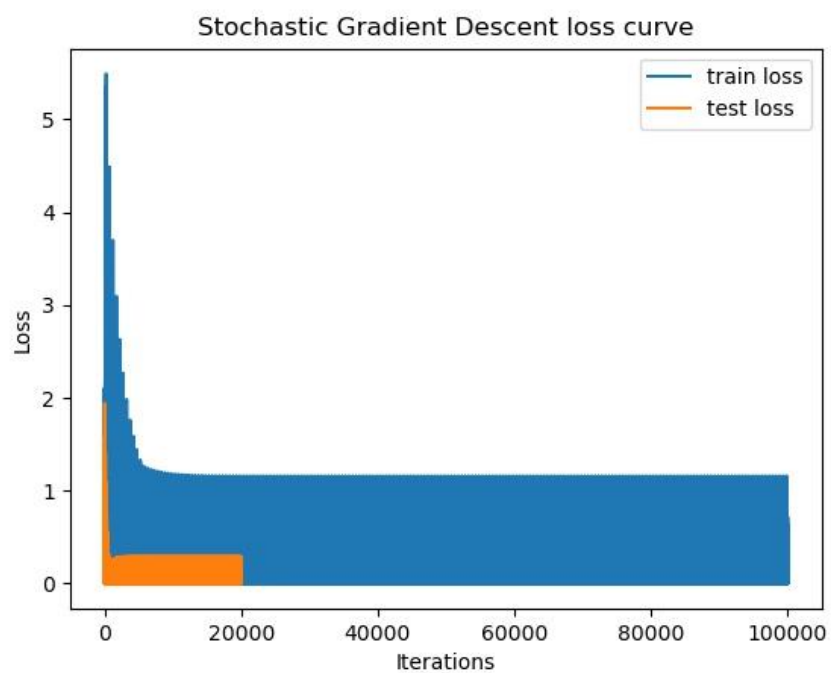
## C. Stochastic Gradient Descent Loss Curve



Figure 3. Stochastic Gradient Descent Loss Curve

## 2. (15%) What's the Mean Square Error of your prediction and ground truth (prediction=model(x_test), ground truth=y_test)

In this section, the situation is divided into three parts, which including Gradient Descent, Minibatch Gradient Descent and Stochastic Gradient Descent. The MSE calculation results of the prediction and the ground truth are shown in Figure 7.

### A. Gradient Descent Mean Square Error

This part of the results can be seen in Figure 4. The training loss in Figure 4 refers to the average of the loss results of each batch when training the model, which means the MSE result between Ground Truth and Prediction. The test loss refers to the average value of each batch of loss results when testing the model, that is, the average value of MSE between the ground truth and the prediction. Because in this method, the batch size will be set to the size of the data, so this MSE between Ground Truth and Prediction will be more accurate. The MSE between Ground Truth and Prediction is 0.00830.

```
100%|                                                    | 1/1 [00:00<00:00, 17476.27it/s]
epoch 997, train loss: 0.13386901619226668, test loss: 0.00830572850313524
100%|                                                    | 1/1 [00:00<00:00, 9709.04it/s]
100%|                                                    | 1/1 [00:00<00:00, 16980.99it/s]
epoch 998, train loss: 0.13386771084797122, test loss: 0.00830720174490206
100%|                                                    | 1/1 [00:00<00:00, 6307.22it/s]
100%|                                                    | 1/1 [00:00<00:00, 14122.24it/s]
epoch 999, train loss: 0.1338664146148669, test loss: 0.008308671067295381
```

Figure 4. Gradient Descent Mean Square Error

## B. Minibatch Gradient Descent Mean Square Error

This part of the results can be seen in Figure 5, and the description is similar to the Figure 4. But in this section, the batch size is 20. The MSE between Ground Truth and Prediction is 0.04124.



```
100%|                                                    | 5/5 [00:00<00:00, 19257.59it/s]
epoch 497, train loss: 0.13441763332272327, test loss: 0.04124741546162789
100%|                                                    | 25/25 [00:00<00:00, 22995.09it/s]
100%|                                                    | 5/5 [00:00<00:00, 29248.98it/s]
epoch 498, train loss: 0.13441763332272327, test loss: 0.04124741546162789
100%|                                                    | 25/25 [00:00<00:00, 23896.44it/s]
100%|                                                    | 5/5 [00:00<00:00, 32413.48it/s]
epoch 499, train loss: 0.13441763332272327, test loss: 0.04124741546162789
```

Figure 5. Minibatch Gradient Descent Mean Square Error

## C. Stochastic Gradient Descent Mean Square Error

This part of the results can be seen in Figure 6, and the description is similar to the Figure 4. The MSE between Ground Truth and Prediction is 0.0399215.



```
100%|                                                    | 100/100 [00:00<00:00, 88301.14it/s]
epoch 197, train loss: 0.13444524231659954, test loss: 0.03992154252589639
100%|                                                    | 500/500 [00:00<00:00, 32069.98it/s]
100%|                                                    | 100/100 [00:00<00:00, 53731.80it/s]
epoch 198, train loss: 0.13444524231659954, test loss: 0.03992154252589639
100%|                                                    | 500/500 [00:00<00:00, 41424.41it/s]
100%|                                                    | 100/100 [00:00<00:00, 105015.12it/s]
epoch 199, train loss: 0.13444524231659954, test loss: 0.03992154252589639
```

Figure 6. Stochastic Gradient Descent Mean Square Error

## D. Calculate MSE between Prediction and Ground truth

The calculation of MSE (x_test and y_test) can be seen in the code of Figure 7. Line 146 is to find the prediction, which is the process of weight*x_test+bias. As

for line 149, it starts to calculate the MSE between prediction and ground truth (y_test). The detailed calculation code of MSE can is shown in Figure 8.

```
142             # testing phase(calculate test loss)
143             for x_batch, y_batch in tqdm(test_batch):
144
145                 # like y_pred = model(x)
146                 y_pred = np.matmul(x_batch, best_theta)
147
148                 # calculate mse between prediction and ground truth
149                 loss = criterion(y_pred, y_batch, batch_size)
150
151                 # save the loss value to plot the figure
152                 test_batch_loss.append(loss)
153                 test_loss_history.append(loss)
```

Figure 7. MSE Calculation

```
14    def criterion(y_hat, y, m):  # calculate the loss through Mean Square Error
15        return 1/(2*m) * np.sum((y_hat - y)**2)
```

Figure 8. MSE Calculation detailed

# 3. (15%) What're the weights and intercepts of your linear model?

In this part, apart from dividing this part into gradient descent, Mini batch gradient descent and stochastic gradient descent, I also divide it into **normalize** and **unnormalize** for discussion.

## A. Gradient Descent weights and intercepts

● Normalize

In this assignment, I have saved the best weight and intercepts according to the MSE, and the best weight in gradient descent method is 4.221022, the intercepts is 0.8018.

```
100%|                                                    | 1/1 [00:00<00:00, 17476.27it/s]
epoch 997, train loss: 0.13386901619226668, test loss: 0.00830572850313524
100%|                                                    | 1/1 [00:00<00:00, 9709.04it/s]
100%|                                                    | 1/1 [00:00<00:00, 16980.99it/s]
epoch 998, train loss: 0.13386771084797122, test loss: 0.00830720174490206
100%|                                                    | 1/1 [00:00<00:00, 6307.22it/s]
100%|                                                    | 1/1 [00:00<00:00, 14122.24it/s]
epoch 999, train loss: 0.1338664146148669, test loss: 0.008308671067295381
The Gradient descent best weight of the linear model is 4.221022743512345 ,and the intercepts is 0.80182619952689
13
```

Figure 9. Gradient Descent weight and intercepts(normalize)

● UnNormalize

The best weight in gradient descent method is 0.8179, the intercepts is 0.7845. In this part, you can see that after normalizing, the weight value will become about 6 times.

```
100%|                                                          | 1/1 [00:00<00:00, 8456.26it/s]
100%|                                                          | 1/1 [00:00<00:00, 18157.16it/s]
epoch 997, train loss: 0.13368200181610046, test loss: 0.00687029732074477
100%|                                                          | 1/1 [00:00<00:00, 10672.53it/s]
100%|                                                          | 1/1 [00:00<00:00, 21076.90it/s]
epoch 998, train loss: 0.13368200181610046, test loss: 0.00687029732074477
100%|                                                          | 1/1 [00:00<00:00, 8848.74it/s]
100%|                                                          | 1/1 [00:00<00:00, 20560.31it/s]
epoch 999, train loss: 0.13368200181610046, test loss: 0.00687029732074477
The Gradient descent best weight of the linear model is 0.8179703770324728 ,and the intercepts is 0.7845650819162
41
```

Figure 10. Gradient Descent weight and intercepts

## B. Minibatch Gradient Descent weights and intercepts

● Normalize

In this assignment, I have saved the best weight and intercepts according to the

MSE, and the best weight in minibatch gradient descent method is 4.321306, the

intercepts is 0.7733.

```
100%|                                                          | 5/5 [00:00<00:00, 19257.59it/s]
epoch 497, train loss: 0.13441763332272327, test loss: 0.04124741546162789
100%|                                                          | 25/25 [00:00<00:00, 22995.09it/s]
100%|                                                          | 5/5 [00:00<00:00, 29248.98it/s]
epoch 498, train loss: 0.13441763332272327, test loss: 0.04124741546162789
100%|                                                          | 25/25 [00:00<00:00, 23896.44it/s]
100%|                                                          | 5/5 [00:00<00:00, 32413.48it/s]
epoch 499, train loss: 0.13441763332272327, test loss: 0.04124741546162789
The Minibatch Gradient Descent best weight of the linear model is 4.321306454481022 ,and the intercepts is 0.7733
880702374907
```

Figure 11. Minibatch Gradient Descent weight and intercepts(normalize)

● UnNormalize

The best weight in gradient descent method is 0.80244, the intercepts is 0.75648.

In this part, you can see that after normalizing, the weight value will become about

6 times.

```
100%|                                                              | 25/25 [00:00<00:00, 29232.67it/s]
100%|                                                              | 5/5 [00:00<00:00, 42711.85it/s]
epoch 497, train loss: 0.13535124205438773, test loss: 0.034816128155235666
100%|                                                              | 25/25 [00:00<00:00, 29208.25it/s]
100%|                                                              | 5/5 [00:00<00:00, 42538.58it/s]
epoch 498, train loss: 0.13535124205438773, test loss: 0.034816128155235666
100%|                                                              | 25/25 [00:00<00:00, 29620.79it/s]
100%|                                                              | 5/5 [00:00<00:00, 43873.47it/s]
epoch 499, train loss: 0.13535124205438773, test loss: 0.034816128155235666
The Minibatch Gradient Descent best weight of the linear model is 0.8024482255801342 ,and the intercepts is 0.756
4840225346317
```

Figure 12. Minibatch Gradient Descent weight and intercepts

## C. Stochastic Gradient Descent weights and intercepts

● Normalize

In this assignment, I have saved the best weight and intercepts according to the MSE, and the best weight in minibatch gradient descent method is 4.3159, the intercepts is 0.755. In Figure 15, we can see that the saving code, weight and intercepts are saved in the "best_theta" parameter.

```
100%|                                                              | 100/100 [00:00<00:00, 88301.14it/s]
epoch 197, train loss: 0.13444524231659954, test loss: 0.03992154252589639
100%|                                                              | 500/500 [00:00<00:00, 32069.98it/s]
100%|                                                              | 100/100 [00:00<00:00, 53731.80it/s]
epoch 198, train loss: 0.13444524231659954, test loss: 0.03992154252589639
100%|                                                              | 500/500 [00:00<00:00, 41424.41it/s]
100%|                                                              | 100/100 [00:00<00:00, 105015.12it/s]
epoch 199, train loss: 0.13444524231659954, test loss: 0.03992154252589639
The Stochastic Gradient Descent best weight of the linear model is 4.315912459202712 ,and the intercepts is 0.755
8936600133507
```

Figure 13. Stochastic Gradient Descent weight and intercepts(normalize)

● UnNormalize

The best weight in gradient descent method is 0.796593, the intercepts is 0.7392. In this part, you can see that after normalizing, the weight value will become about 6 times.

```
100%|                                                            | 500/500 [00:00<00:00, 49750.95it/s]
100%|                                                            | 100/100 [00:00<00:00, 99391.09it/s]
epoch 197, train loss: 0.1350201054720081, test loss: 0.035272975863546344
100%|                                                            | 500/500 [00:00<00:00, 48740.37it/s]
100%|                                                            | 100/100 [00:00<00:00, 103768.04it/s]
epoch 198, train loss: 0.1350201054720081, test loss: 0.035272975863546344
100%|                                                            | 500/500 [00:00<00:00, 54989.96it/s]
100%|                                                            | 100/100 [00:00<00:00, 112508.15it/s]
epoch 199, train loss: 0.1350201054720081, test loss: 0.035272975863546344
The Stochastic Gradient Descent best weight of the linear model is 0.7965938046686774 ,and the intercepts is 0.73
92428834352359
```

Figure 14. Stochastic Gradient Descent weight and intercepts

```python
117         for epoch in range(epochs):
118
119             train_batch_loss, test_batch_loss = [], []
120
121             # training phase(calculate train loss)
122             for x_batch, y_batch in tqdm(train_batch):
123
124                 # something like y_pred = model(x)
125                 y_pred = np.matmul(x_batch, theta)
126
127                 # calculate mse between prediction and ground truth
128                 loss = criterion(y_pred, y_batch, batch_size)
129
130                 # calculate the gradient and do backpropagation
131                 theta -= lr * \
132                     compute_gradient(x_batch, y_pred, y_batch, batch_size)
133
134                 # save the loss value to plot the figure
135                 train_batch_loss.append(loss)
136                 train_loss_history.append(loss)
137
138                 # save the best weight and bias
139                 if min_loss > loss:
140                     best_theta = theta
```

Figure 15. Saving the best weights and intercepts

## 4. (10%) What's the difference between Gradient Descent, Mini-Batch Gradient Descent, and Stochastic Gradient Descent?

The difference between the three methods is the batch size. The batch size used by gradient descent is the number of training data (500), the batch size of mini-

batch gradient descent is selected from 1 to 500, and the batch size of stochastic

gradient descent is 1. The setting of the hyperparameters including learning rate(lr),

batch size(batch_size) ,and number of epochs(epochs) can be seen in the Figure

22.

## A. Gradient Descent weights and intercepts

The figure below is about data point and fitting line, especially during training

and testing. During the training stage, because the batch size of the gradient descent

is an integer of the data, there should be more epochs during the training.



Figure 16. Training data point and fitting line

Figure 17. Testing data point and fitting line

## B. Minibatch Gradient Descent weights and intercepts

The figure below is about data point and fitting line, especially during training and testing. During the training stage, because the batch size of the Minibatch Gradient Descent is the number during 1~whole data, there can have less epochs during the training.

Figure 18. Training data point and fitting line



Figure 19. Testing data point and fitting line

## C. Stochastic Gradient Descent weights and intercepts

The figure below is about data point and fitting line, especially during training and testing. During the training stage, because the batch size of the Stochastic Gradient Descent is 1, the number of epochs can less than minibatch gradient descent.



Figure 20. Training data point and fitting line

Figure 21. Testing data point and fitting line

```
72        # Gradient descent mode
73        if mode == 1:
74            lr = 1e-1
75            epochs = 1000
76            batch_size = 500
77            name = "Gradient descent"
78
79        # Minibatch Gradient Descent mode
80        elif mode == 2:
81            lr = 1e-1
82            epochs = 500
83            batch_size = 20
84            name = "Minibatch Gradient Descent"
85
86        # Stochastic Gradient Descent mode
87        elif mode == 3:
88            lr = 1e-2
89            epochs = 200
90            batch_size = 1
91            name = "Stochastic Gradient Descent"
```

Figure 22. The hyperparameters in each gradient descent method

# ● Part. 2, Questions (40%)

系所:資料工所.    學號:31055103|    姓名:張皓雯.

Part 2. Questions.

1.

|  | apples | oranges | guavas. | 選到)機率 |
|---|---|---|---|---|
| R. | 3 | 4 | 3. | 0.2. |
| B | 2 | 0 | 2. | 0.4. |
| G. | 12 | 4 | 4. | 0.4. |

(1). guava.          Red          Blue          Green.

$$P(g) = P(R)P(g|R) + P(B)P(g|B) + P(G)P(g|G)$$

$$= 0.2 \times 0.3 + 0.4 \times 0.5 + 0.4 \times 0.2$$

$$= 0.34. \text{ ※}$$

(2).  Blue  apple

$$P(B|a) = \frac{P(a|B)P(B)}{P(a)} = \frac{P(a|B)P(B)}{P(a|R)P(R) + P(a|B)P(B) + P(a|G)P(G)}$$

$$= \frac{0.5 \times 0.4}{0.3 \times 0.2 + 0.5 \times 0.4 + 0.6 \times 0.4}$$

$$= \frac{0.2}{0.5}$$

$$= 0.4 \text{ ※}$$

2.

(2) $P(x, C_1) \geq P(x, C_2)$    |    $x \in R_1$

$P(x, C_2) \geq P(x, C_1)$    |    $x \in R_2$

$\int \{ P(x, C_1) P(x, C_2) \}^{\frac{1}{2}} dx$

$= \int_{R_1} \{ P(x, C_1) P(x, C_2) \}^{\frac{1}{2}} dx + \int_{R_2} \{ P(x, C_1) P(x, C_2) \}^{\frac{1}{2}} dx \geq \int_{R_1} P(x, C_2) dx + \int_{R_2} P(x, C_1)$

$\phantom{=========================}\parallel$

$\phantom{==============================}P(\text{mistake})$

$\phantom{================================}\#$

(1).    prove.

if $a \leq b$, then $a \leq (ab)^{\frac{1}{2}}$

$(ab)^{\frac{1}{2}} \geq (a \times a)^{\frac{1}{2}} = a$

$\therefore \sqrt{ab} \geq a$.

---

3. (1) prove $E(x) = E_y [E_x [x | y]]$.

$= E_y [ \sum_x x \cdot P(x = x | Y = y) ]$

$= \sum_y \sum_x x \cdot P(X = x | Y = y) P(Y = y)$

$= \sum_y \sum_x x \cdot P(Y = y | X = x) P(X = x)$

$= \sum_x x \cdot P(X = x) \sum_y P(Y = y | X = x)$

$= \sum_x x P(X = x)$

$= E(x)$.

3.12. prove $\quad Var[x] = E_y[Var_x[x|y]] + Var_y[E_x[x|y]]$

$\text{H}'\lambda \cdot Var(x) = E[x^2] - [E[x]]^2$

$E_y[Var_x[x|y]] + Var_y(E_x[x|y])$

$= E_y[E_x[x^2|y] - (E_x[x|y])^2] + E_y[(E_x[x|y])^2] - (E_y[E_x[x|y]])^2$

$= E_x[x^2] - E_y[[E_x[x|y]]^2] + E_y[(E_x[x|y])^2] - (E_x(x))^2$

$= E_x[x^2] - (E_x[x])^2$

$= Var_x[x]$

# ● Part. 3, Model Architecture and Method

## 1. Flow chart

```python
29  def data_prepare(csv_name, norm):  # pack input into several batch (dataloader)
30
31      x_data, y_data = [], []
32
33      # loading training data into several batch
34      if "train" in csv_name:
35          df_train = pd.read_csv(csv_name)
36          x, y = df_train['x_train'].to_numpy(), df_train['y_train'].to_numpy()
37
38      # loading testing data into several batch
39      elif "test" in csv_name:
40          df_test = pd.read_csv(csv_name)
41          x, y = df_test['x_test'].to_numpy(), df_test['y_test'].to_numpy()
42
43      plot_x = x
44      plot_y = y
45
46      # feature scaling
47      if norm == "y":
48          x = normalize(x)
49
50      x = x.reshape(-1, 1)
51      y = y.reshape(-1, 1)
52
53      x = np.hstack((np.ones(x.shape), x))
54
55      for curr_batch in range(0, len(x), batch_size):
56          x_data.append(x[curr_batch:curr_batch+batch_size])
57          y_data.append(y[curr_batch:curr_batch+batch_size])
58
59      x_data = np.array(x_data)
60      y_data = np.array(y_data)
61
62      batch = list(zip(x_data, y_data))
63
64      return plot_x, plot_y, batch
```

```python
9   # scale feature into [0,1] through Mean Normalization (accelerate diverge)
10  def normalize(x):
11      return (x - np.mean(x, axis=0))/(np.max(x, axis=0) - np.min(x, axis=0))
```

```python
109         # random the initial weight and bias(storage weight)
110         theta = np.random.randn(2, 1)
```

```python
72      # Gradient descent mode
73      if mode == 1:
74          lr = 1e-1
75          epochs = 1000
76          batch_size = 500
77          name = "Gradient descent"
78
79      # Minibatch Gradient Descent mode
80      elif mode == 2:
81          lr = 1e-1
82          epochs = 500
83          batch_size = 20
84          name = "Minibatch Gradient Descent"
85
86      # Stochastic Gradient Descent mode
87      elif mode == 3:
88          lr = 1e-2
89          epochs = 200
90          batch_size = 1
91          name = "Stochastic Gradient Descent"
```

```python
117     for epoch in range(epochs):
118
119         train_batch_loss, test_batch_loss = [], []
120
121         # training phase(calculate train loss)
122         for x_batch, y_batch in tqdm(train_batch):
123
124             # something like y_pred = model(x)
125             y_pred = np.matmul(x_batch, theta)
126
127             # calculate mse between prediction and ground truth
128             loss = criterion(y_pred, y_batch, batch_size)
129
130             # calculate the gradient and do backpropagation
131             theta -= lr * \
132                 compute_gradient(x_batch, y_pred, y_batch, batch_size)
133
134             # save the loss value to plot the figure
135             train_batch_loss.append(loss)
136             train_loss_history.append(loss)
137
138             # save the best weight and bias
139             if min_loss > loss:
140                 best_theta = theta
```

Flow chart (boxes):
- Start
- Choose Gradient Descent method
- Hyperparameters setting
- Pack training and testing data into several batch
- Check if you want to normalize the input — True / False
- Scaling the feature into [0,1] (Normalize)
- initial weight and bias
- (line 125) Calculate the prediction with weight and bias
- (line 128) Calculate MSE
- (line 131) Calculate gradient and update the weight, bias
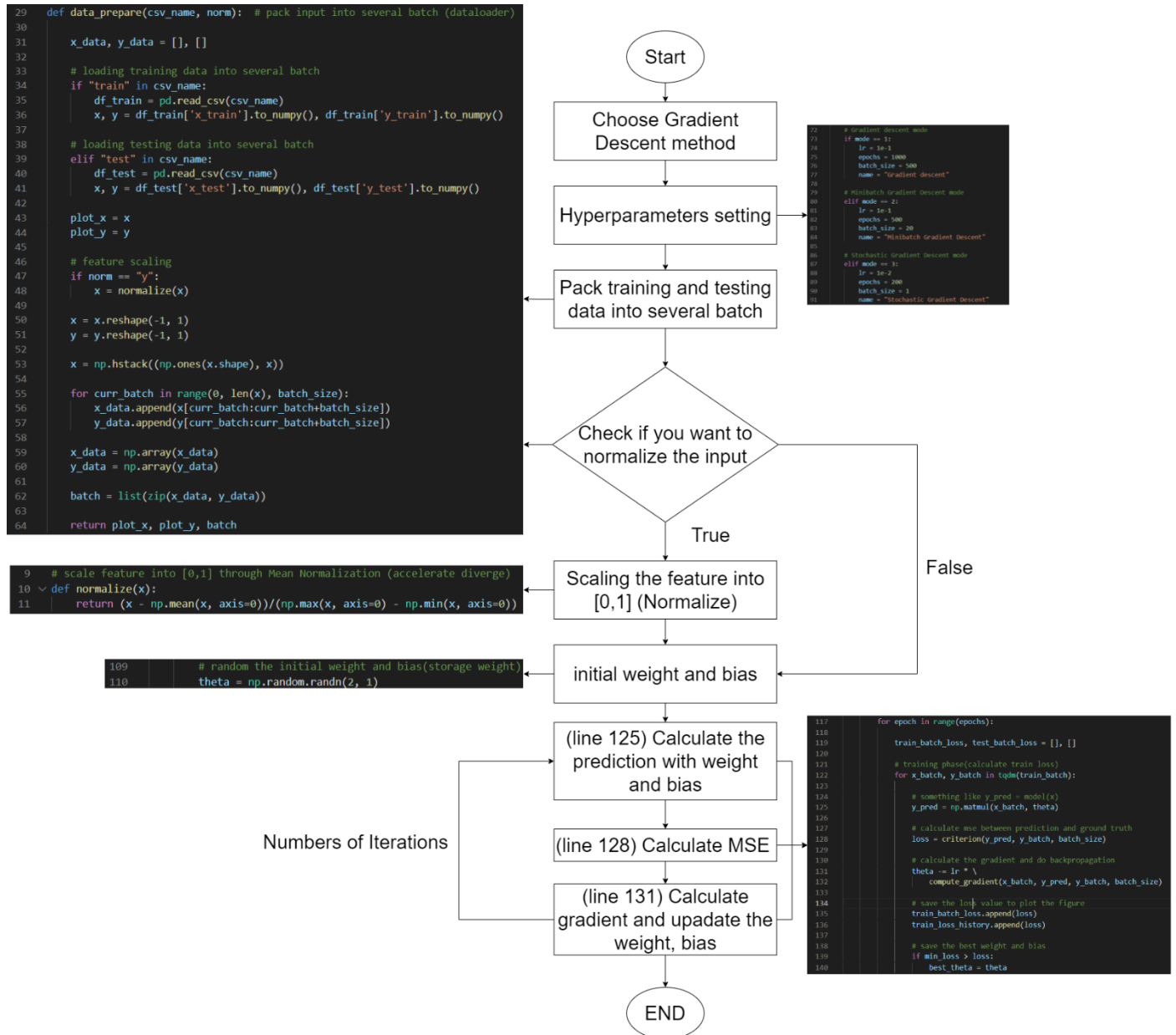- Numbers of Iterations
- END

Figure 23. Programming flow chart

## 2. Normaliztion(feature scaling)

Normalization is to subtract the average of all x-pointoordinates from each x -

point coordinate and then divide it by the interval between the maximum and

minimum values. If use this method, the model will converge more quickly.

```
 9     # scale feature into [0,1] through Mean Normalization (accelerate diverge)
10   ∨ def normalize(x):
11         return (x - np.mean(x, axis=0))/(np.max(x, axis=0) - np.min(x, axis=0))
```

Figure 24. Normalization method