

# AI and Machine Learning, Homework2

Author: Ying Yiwen  
Number: 12210159

## Contents

<b>1</b>	<b>Code Implementation</b>	<b>2</b>
1.1	MSE Loss . . . . .	2
1.2	Update Method . . . . .	2
1.3	Normalization . . . . .	3
1.4	Random Split . . . . .	4
1.5	Main Function . . . . .	4
<b>2</b>	<b>Analyzation of Different Methods</b>	<b>5</b>
2.1	Different Update Methods . . . . .	5
2.2	Different Normalization Method . . . . .	6

# 1 Code Implementation

## 1.1 MSE Loss

use the formula:

$$l(\omega) = \frac{(y_{pred} - y_{true})^2}{2N} \quad (1)$$

```
1 def calculate_loss(self, y, y_pred):
2     # MSE loss
3     N = y.size
4     return np.sum((y_pred - y) ** 2) / (2 * N)
```

## 1.2 Update Method

train function with BGD method:

```
1 def train_BGD(self, X_train, y_train, X_test, y_test):
2     # BGD
3     self.W = 10*np.random.rand(X_train.shape[1] + 1)
4     X_train, X_test = self.preprocess_data_X(X_train), self.preprocess_data_X(
5         X_test)
6     for _ in range(self.n_iter):
7         # compute loss
8         y_pred = self.predict(X_train)
9         train_loss = self.calculate_loss(y_train, y_pred)
10        self.train_loss.append(train_loss)
11        y_pred_test = self.predict(X_test)
12        test_loss = self.calculate_loss(y_test, y_pred_test)
13        self.test_loss.append(test_loss)
14        # gradient descent
15        self.W -= self.lr * self.gradient(X_train, y_train, y_pred)
```

train function with SGD method:

```
1 def train_SGD(self, X_train, y_train, X_test, y_test):
2     # SGD
3     self.W = 10*np.random.rand(X_train.shape[1] + 1)
4     X_train, X_test = self.preprocess_data_X(X_train), self.preprocess_data_X(
5         X_test)
6     N = y_train.size
7     for _ in range(self.n_iter):
8         # shuffle the data
9         indices = np.random.permutation(N)
10        for i in indices:
11            x_i = X_train[i:i+1]
12            y_i = y_train[i:i+1]
13            # compute loss
14            y_pred = self.predict(x_i)
15            train_loss = self.calculate_loss(y_i, y_pred)
16            self.train_loss.append(train_loss)
17            y_pred_test = self.predict(X_test)
18            test_loss = self.calculate_loss(y_test, y_pred_test)
19            self.test_loss.append(test_loss)
20        # compute gradient
```

```

20         grad = self.gradient(x_i, y_i, y_pred)
21         self.W -= self.lr * grad

```

train function with MBGD method:

```

1 def train_MBGD(self, X_train, y_train, X_test, y_test):
2     # MBGD
3     self.W = 10*np.random.rand(X_train.shape[1] + 1)
4     X_train, X_test = self.preprocess_data_X(X_train), self.preprocess_data_X(
        X_test)
5     N = y_train.size
6     for _ in range(self.n_iter):
7         # shuffle the data
8         indices = np.random.permutation(N)
9         for start in range(0, N, self.batch_size):
10             end = min(start + self.batch_size, N) # not out of range
11             # get batch data
12             batch_indices = indices[start:end]
13             X_batch = X_train[batch_indices]
14             y_batch = y_train[batch_indices]
15             # predict and compute loss
16             y_pred = self.predict(X_batch)
17             train_loss = self.calculate_loss(y_batch, y_pred)
18             self.train_loss.append(train_loss)
19             y_pred_test = self.predict(X_test)
20             test_loss = self.calculate_loss(y_test, y_pred_test)
21             self.test_loss.append(test_loss)
22             # compute gradient
23             grad = self.gradient(X_batch, y_batch, y_pred)
24             self.W -= self.lr * grad

```

### 1.3 Normalization

There should be normalization function for x and inverse normalization for W. As follows:  
min-max normalization:

```

1 def min_max_normalization(self, x):
2     # normalize the data by min-max normalization
3     _min = np.min(x, axis=0)
4     _max = np.max(x, axis=0)
5     _range = _max - _min
6     normalized = (x - _min) / _range
7     return normalized, _min, _max
8
9 def inverse_min_max_weight(self, W, _min, _max):
10    # inverse min-max normalization
11    W[:1] -= W[1:] * _min / (_max - _min)
12    W[1:] /= (_max - _min)
13    return W

```

mean normalization:

```

1 def mean_normalization(self, x):
2     # normalize the data by mean normalization
3     mu = np.mean(x, axis=0)

```

```

4     sigma = np.std(x, axis=0)
5     normalized = (x - mu) / sigma
6     return normalized, mu, sigma
7
8 def inverse_mean_weight(self, W, mu, sigma):
9     # inverse mean normalization
10    W[:1] -= W[1:] * mu / sigma
11    W[1:] /= sigma
12    return W

```

## 1.4 Random Split

The input data should be split into train set and test set randomly. Shuffle the data also prevent the fitting pay importance to too few features.

```

1 def random_split(self, x, y, test_size=0.1, random_state=None):
2     # shuffle the data
3     if random_state is not None:
4         np.random.seed(random_state)
5         indices = np.arange(len(x))
6         np.random.shuffle(indices)
7     # split the data
8     split_index = int(len(x) * (1-test_size))
9     train_indices, test_indices = indices[:split_index], indices[split_index:]
10    xtrain, ytrain = x[train_indices], y[train_indices]
11    xtest, ytest = x[test_indices], y[test_indices]
12    return xtrain, ytrain, xtest, ytest

```

## 1.5 Main Function

With all the function written, here's the main function:

```

1 if __name__ == "__main__":
2     # record the time
3     start_time = time.time()
4
5     # generate data
6     X = np.arange(100).reshape(100,1)
7     a, b = 1, 10
8     y = a * X + b + np.random.normal(0, 5, size=X.shape)
9     y = y.reshape(-1)
10
11    # set parameters
12    n_iter, lr = 30000, 5e-4
13    model = LinearRegression(n_iter=n_iter, lr=lr)
14
15    # normalize the data
16    method = 'min_max'
17    X_normalized, param1, param2 = model.normalize(X, method=method)
18
19    # split the data
20    X_train, y_train, X_test, y_test = model.random_split(X_normalized, y,
        test_size=0.2)

```

```

21
22 # train the model
23 model.train(X_train, y_train, X_test, y_test, method='MBGD')
24 model.plot_loss()
25 end_time = time.time()
26
27 # compute loss
28 train_loss = model.calculate_loss(y_train, model.predict(model.
    preprocess_data_X(X_train)))
29 test_loss = model.calculate_loss(y_test, model.predict(model.preprocess_data_X
    (X_test)))
30
31 # inverse normalization
32 if method == 'min_max':
33     W_original = model.inverse_min_max_weight(model.W, param1, param2)
34 elif method == 'mean':
35     W_original = model.inverse_mean_weight(model.W, param1, param2)
36 else:
37     W_original = model.W
38
39 # output the result
40 model.plot_fit(X, y)
41 print(f'n of iteration: {n_iter}, loss rate: {lr}')
42 print(f'Learned weights: {W_original}, Training loss: {train_loss}, Testing
    loss: {test_loss}')
43 print(f'Time: {end_time - start_time}s')

```

## 2 Analyzation of Different Methods

### 2.1 Different Update Methods

To control variable, we don't use any normalization method, to see the difference of update methods only.

MBGD:

n of iteration: 10000, loss rate: 0.0002

Learned weights: [10.15597587 1.0220182 ], Training loss: 14.307060377519765, Testing loss: 9.575175286124201

Time: 4.011836767196655s

BGD:

n of iteration: 50000, loss rate: 0.0001

Learned weights: [9.80610593 1.00979871], Training loss: 11.736946331559594, Testing loss: 13.800069762104055

Time: 3.438744068145752s

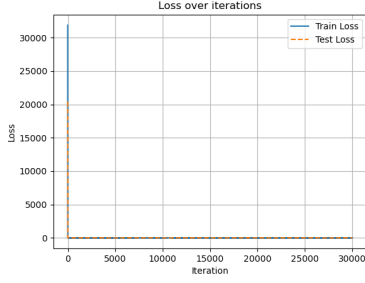
SGD:

n of iteration: 2000, loss rate: 0.0001

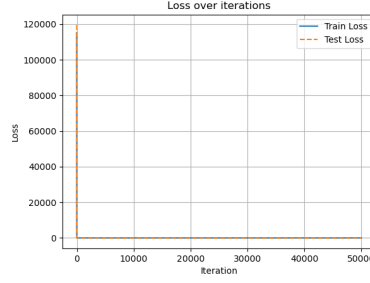
Learned weights: [10.24699952 1.03413564], Training loss: 12.558711291049434, Testing loss: 8.769263129234005

Time: 6.4901885986328125s

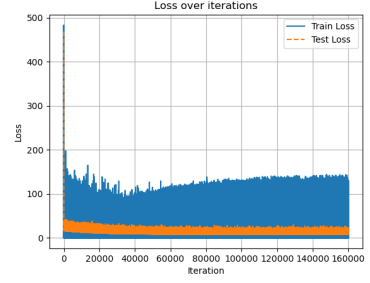
BGD can compute parallelly, but it use too much samples at a time, so it works slowly and needs large memory consume. BGD also may falls to local minimum as the loss function may be to smooth. But BGD can cope with input noises, so it can converge well.



(a) MBGD Loss Function

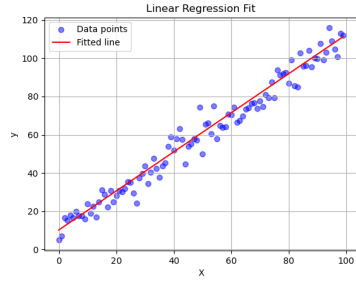


(b) BGD Loss Function

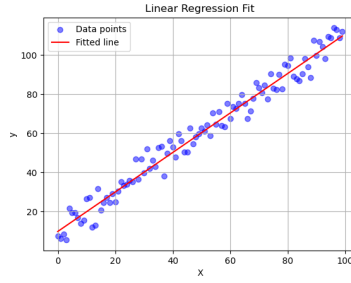


(c) SGD Loss Function

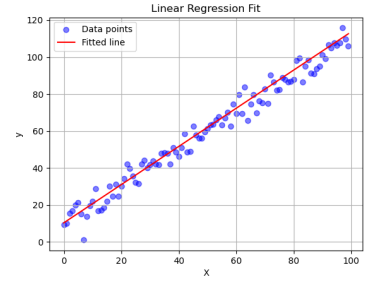
Fig. 1: The Trend of Loss Function of Different Gradient Method



(a) MBGD Fitting Results



(b) BGD Fitting Results



(c) SGD Loss Function

Fig. 2: The Fitting Results of Different Gradient Method

SGD only use one sample at a time, so it uses less memory and works quickly. Thanks to the noises in input data, SGD can easily jump out of local minimum. However, when it actually converges, those noises may make SGD harder to get the global minimum, so it doesn't converge very well.

MBGD take the advantage of both methods, but it need suitable choice of batch size to balance the influence.

## 2.2 Different Normalization Method

To control variable, we all use MBGD method in this section to only compare the difference of normalization method.

No Normalization Method:

n of iteration: 50000, loss rate: 0.0002

Learned weights: [9.0494355 1.02331261], Training loss: 10.862699339670602, Testing loss: 17.802074655015762

Time: 6.295790195465088s

Min-Max:

n of iteration: 30000, loss rate: 0.0005

Learned weights: [10.72168941 0.97626939], Training loss: 12.24969162788654, Testing loss: 31.357194692974275

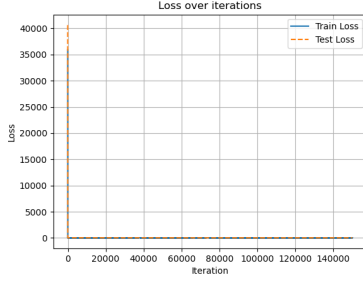
Time: 4.983108997344971s

Mean:

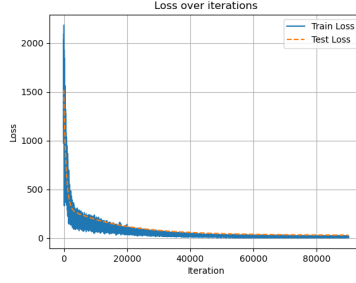
n of iteration: 10000, loss rate: 0.0004

Learned weights: [10.84252905 1.00589178], Training loss: 13.687261666680957, Testing loss: 14.04630888716398

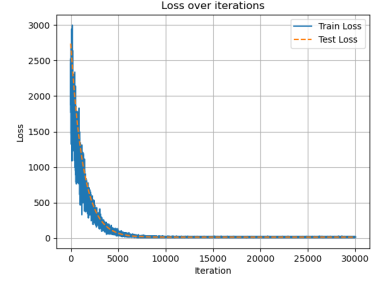
Time: 6.000371217727661s



(a) No Normalization

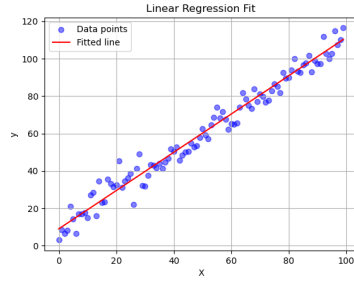


(b) Min-Max Normalization

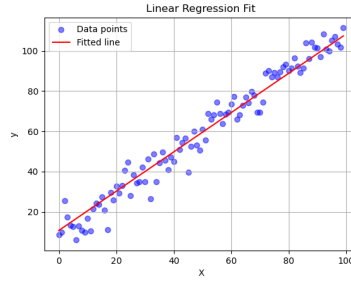


(c) Mean Normalization

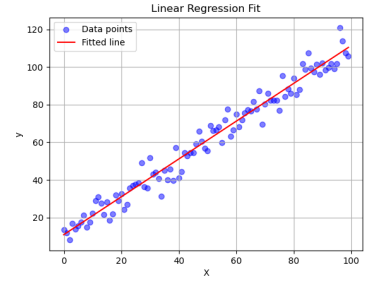
Fig. 3: The Trend of Loss Function of Different Gradient Method



(a) No Normalization



(b) Min-Max Normalization



(c) Mean Normalization

Fig. 4: The Fitting Results of Different Gradient Method

Normalization actually changes the step of each element of  $W$ . When normalization can balance the influence of each feature by correct choice of normalization, it can work well on inputs of different order of magnitude, and thus, it makes fitting better and quicker. However, it also lose efficacy when there is extreme outlier. Normalization also needs inverse normalization after fitting, which makes the algorithm more complex.

Min-Max Normalization can eliminate the effects of magnitude and order of magnitude. However, it may make the method overly dependent on two extreme values when changing the weights of the variables.

Mean Normalization can remove the quantisation. Dimensionlessness also eliminates the differences in the degree of variation of each variable, so that the converted variables are treated equally in the cluster analysis. However, in actual analyses, the importance of each variable in the analysis shouldn't be exactly the same. At the same time, mean normalization only work well with gaussian input, otherwise it may makes the matter rather worse.