

Project1 Report

Xueying Sui

1001682442

Problem

We need to train a linear regression model, and use this model to make prediction. We have a data set, this data set contains 150 data items. We should divide it into two parts: the training set and the test set.

We use the training data train our model, then we can get a set of parameters of estimate function. We use this estimate function and test data to make prediction and calculate accuracy.

Data

We have 150 data items, these items have three tags: Iris-setosa, Iris-versicolor and Iris-virginica. Each of the tags corresponds to 50 data items. To train linear regression model, we should assign values to these tags, like (0, 1, 2) or (100, 101, 102), etc.

When we dividing the data set, we can try five or five separate or four or six separate.

Method 1

$$\theta = (A^T A)^{-1} A^T Y$$

We can see the input data points of training data as a matrix A, see the label data as a matrix Y, and see parameters as a matrix θ . So the sum of the distances between the model and each point of the training data can be expressed as:

$$\begin{aligned} distance &= (A\theta - Y)^T (A\theta - Y) \\ &= (\theta^T A^T - Y^T) (A\theta - Y) \\ &= \theta^T A^T A\theta - \theta^T A^T Y - Y^T A\theta + Y^T Y \\ &= \theta^T A^T A\theta - 2Y^T A\theta + Y^T Y \end{aligned}$$

Now our goal is to get the minimum value of distance:

$$\frac{\partial distance}{\partial \theta} = 2\theta^T A^T A - 2Y^T A$$

We make the result of above function equals 0, we can get:

$$\begin{aligned} \theta^T A^T A &= Y^T A \Rightarrow A^T A\theta = A^T Y \Rightarrow (A^T A)^{-1} A^T A\theta = (A^T A)^{-1} A^T Y \\ \Rightarrow \theta &= (A^T A)^{-1} A^T Y \end{aligned}$$

Now we can directly use this function to calculate parameters. And we can use these parameters get estimate function, like this:

$$y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2 + \theta_3 * x_3 + \theta_4 * x_4$$

For more convenient matrix calculation, we can make $x_0=1$, and the estimate function like this:

$$y = \theta_0 * x_0 + \theta_1 * x_1 + \theta_2 * x_2 + \theta_3 * x_3 + \theta_4 * x_4$$

(1) We use 90 data items as training data(including 3 groups, each group has 30 data items), 60 data item as testing data(including 3 groups, each group has 20 data items).

(i) We assign values to the tags. Replace Iris-setosa with 0, replace Iris-versicolor with 1 and replace Iris-virginica with 2

Exp: for one of the data items (5.1,3.5,1.4,0.2,Iris-setosa), we can write like this:

$$\theta_0 * 1 + \theta_1 * 5.1 + \theta_2 * 3.5 + \theta_3 * 1.4 + \theta_4 * 0.2 = 0$$

for one of the data items (7.0,3.2,4.7,1.4,Iris-versicolor), we can write like this:

$$\theta_0 * 1 + \theta_1 * 7.0 + \theta_2 * 3.2 + \theta_3 * 4.7 + \theta_4 * 1.4 = 1$$

for one of the data items (6.3,3.3,6.0,2.5,Iris-virginica), we can write like this:

$$\theta_0 * 1 + \theta_1 * 6.3 + \theta_2 * 3.3 + \theta_3 * 6.0 + \theta_4 * 2.5 = 2$$

We use training data calculate the values of these parameters. After we get the values of parameters, we can use testing data to test the accuracy of this linear regression model.

Firstly, I use MSE(Mean Squared Error) to calculate the loss of training data and estimate the loss of testing data.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Then I use testing data to calculate the accuracy, the method is to separately calculate distances between the estimated value of y and the three tag values, the tag with the smallest distance is used as the estimated tag of this data item, then compared with the real tag. If the estimated tag is the same as real tag, the prediction is correct, if not, the prediction is wrong.

Results

Using the above method we can get the following results.

```
The prediction values of theta are: [[ 0.40679793 -0.17208526 -0.0295757  0.26246741  0.58187472]]
Estimate function is: y = 0.40679793479282894 + ( -0.17208525807795388 ) * x1 + ( -0.029575698828359823 ) * x2 + 0.2624674094814906 * x3 + 0.5818747209357259 * x4
The loss of the estimation function to the training set is  0.048169565780434284
The loss of the estimation function to the testing set is  0.0450385513930922
Accuracy rate is: 0.9666666666666667
```

The result of the loss for training data is about 0.048, and the result of the loss for testing data is about 0.045. These two loss results are similar.

Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%.

(ii) If we exchange training data and testing data, we can get the following results:

```
The prediction values of theta are: [[-0.25161073 -0.00333072 -0.05802762  0.21357679  0.55244971]]
Estimate function is: y = -0.25161073330041184 + ( -0.0033307217412370527 ) * x1 + ( -0.05802761536343806 ) * x2 + 0.2135767870079982 * x3 + 0.5524497111783413 * x4
The loss of the estimation function to the training set is  0.04120457485939926
The loss of the estimation function to the testing set is  0.05383420175538002
Accuracy rate is: 0.9555555555555556
```

Different from before results, the result of the loss for this time's training data is about 0.041, and the result of the loss for this time's testing data is about 0.054.

Then I use the method mentioned above to calculate the accuracy, the result is about 95.6%.

(iii) We assign values to the tags, Replace Iris-setosa with 100, replace Iris-versicolor with 101 and replace Iris-virginica with 102

Similar to (i), we changed the values of the tags, we got the following results:

```
The prediction values of theta are: [[ 1.00406798e+02 -1.72085258e-01 -2.95756988e-02  2.62467409e-01  5.81874721e-01]]
Estimate function is: y = 100.40679793478543 + ( -0.1720852580761849 ) * x1 + ( -0.02957569882829958 ) * x2 + 0.2624674094809283 * x3 + 0.581874720937055 * x4
The loss of the estimation function to the training set is  0.0481695657804343
The loss of the estimation function to the testing set is  0.045038551392913954
Accuracy rate is: 0.9666666666666667
```

The result of the loss for training data is about 0.048, and the result of the loss for testing data is about 0.045. These two loss results are similar.

Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%.

(iv) If we exchange training data and testing data, we can get the following results:

```

The prediction values of theta are: [[ 9.97483893e+01 -3.33072173e-03 -5.80276154e-02  2.13576787e-01
 5.52449711e-01]]
Estimate function is: y = 99.74838926667435 + ( -0.0033307217348408358 ) * x1 + ( -0.05802761536442347 ) * x2 + 0.21357678700595528 * x3 + 0.5524497111783901 * x4
The loss of the estimation function to the training set is  0.04120457485939296
The loss of the estimation function to the testing set is  0.053834201755896154
Accuracy rate is: 0.9555555555555556

```

Different from (iii) results, but similar to (ii) results, the result of the loss for this time's training data is about 0.041, and the result of the loss for this time's testing data is about 0.054.

Then I use the method mentioned above to calculate the accuracy, the result is about 95.6%.

(2) We use 75 data items as training data, 75 data item as testing data(just cut from middle).

(i) We assign values to the tags. Replace Iris-setosa with 0, replace Iris-versicolor with 1 and replace Iris-virginica with 2

Exp: for one of the data items (5.1,3.5,1.4,0.2,Iris-setosa), we can write like this:

$$\theta_0 * 1 + \theta_1 * 5.1 + \theta_2 * 3.5 + \theta_3 * 1.4 + \theta_4 * 0.2 = 0$$

for one of the data items (7.0,3.2,4.7,1.4,Iris-versicolor), we can write like this:

$$\theta_0 * 1 + \theta_1 * 7.0 + \theta_2 * 3.2 + \theta_3 * 4.7 + \theta_4 * 1.4 = 1$$

for one of the data items (6.3,3.3,6.0,2.5,Iris-virginica), we can write like this:

$$\theta_0 * 1 + \theta_1 * 6.3 + \theta_2 * 3.3 + \theta_3 * 6.0 + \theta_4 * 2.5 = 2$$

But thus time, our training data have 75 data items, including all data items that tag equals 0 and a half of data items that tag equals 1, our testing data also have 75 data items, including a half of data items that tag equals 1 and all data items that tag equals 2.

We use training data calculate the values of these parameters. After we get the values of parameters, we can use testing data to test the accuracy of this linear regression model.

Firstly, I use MSE(Mean Squared Error) to calculate the loss of training data and estimate the loss of testing data.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Then I use testing data to calculate the accuracy, the method is to separately calculate distances between the estimated value of y and the three tag values, the tag with the smallest distance is used as the estimated tag of this data item, then compared with the real tag. If the estimated tag is the same as real tag, the prediction is correct, if not, the prediction is wrong.

Results

Using the above method we can get the following results.

```

The prediction values of theta are: [[ 0.10515952 -0.01478308 -0.12074581  0.22881839  0.23218836]]
Estimate function is: y = 0.10515952377528776 + ( -0.014783082799356306 ) * x1 + ( -0.12074580612688214 ) * x2 + 0.2288183933587042 * x3 + 0.23218835590915554 * x4
The loss of the estimation function to the training set is  0.006726888439342197
The loss of the estimation function to the testing set is  0.2666010859890495
Accuracy rate is: 0.3333333333333333

```

The result of the loss for training data is about 0.007, and the result of the loss for testing data is about 0.267. These two loss results are similar.

Then I use the method mentioned above to calculate the accuracy, the result is about 33.3%.

(ii) If we exchange training data and testing data, we can get the following results:

```

The prediction values of theta are: [[ 0.70439344 -0.1962218 -0.29830208  0.37988637  0.63509057]]
Estimate function is: y = 0.7043934414280209 + ( -0.1962218016307008 ) * x1 + ( -0.2983020837257271 ) * x2 + 0.3798863745648632 * x3 + 0.6350905741863688 * x4
The loss of the estimation function to the training set is  0.05309568655301076
The loss of the estimation function to the testing set is  0.27079484173324414
Accuracy rate is: 0.30666666666666664

```

The result of the loss for this time's training data is about 0.053, and the result of the loss for this time's testing data is about 0.271.

Then I use the method mentioned above to calculate the accuracy, the result is about 30.67%.

Method 2

Gradient Descent

First we can make $x_0=1$ and define a cost function:

$$J(\theta_0, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

We use the following formula to update the values of parameters:

$$\theta_i = \theta_i - \alpha \frac{\partial J}{\partial \theta_i}$$

In this formula, α is learning rate. So next we set the value of the learning rate, give parameters(Θ) a random set of initial values, and then set the number of iterations.

Then we need to process the data, we want most x values to be between -1 and 1, so we make standardization use the following formula:

$$x_i = \frac{x_i - \mu}{\sigma}$$

We use 90 data items as training data(including 3 groups, each group has 30 data items), 60 data item as testing data(including 3 groups, each group has 20 data items).

(i) We assign values to the tags. Replace Iris-setosa with 0, replace Iris-versicolor with 1 and replace Iris-virginica with 2

We used the above mentioned method to update the values of parameters, after we get the values of parameters, we can use testing data to test the accuracy of this linear regression model.

Firstly, I used cost function to calculate the loss of training data and estimate the loss of testing data.

Then I use testing data to calculate the accuracy, the method is to separately calculate distances between the estimated value of y and the three tag values, the tag with the smallest distance is used as the estimated tag of this data item, then compared with the real tag. If the estimated tag is the same as real tag, the prediction is correct, if not, the prediction is wrong.

Results

Using the above method we can get the following results.

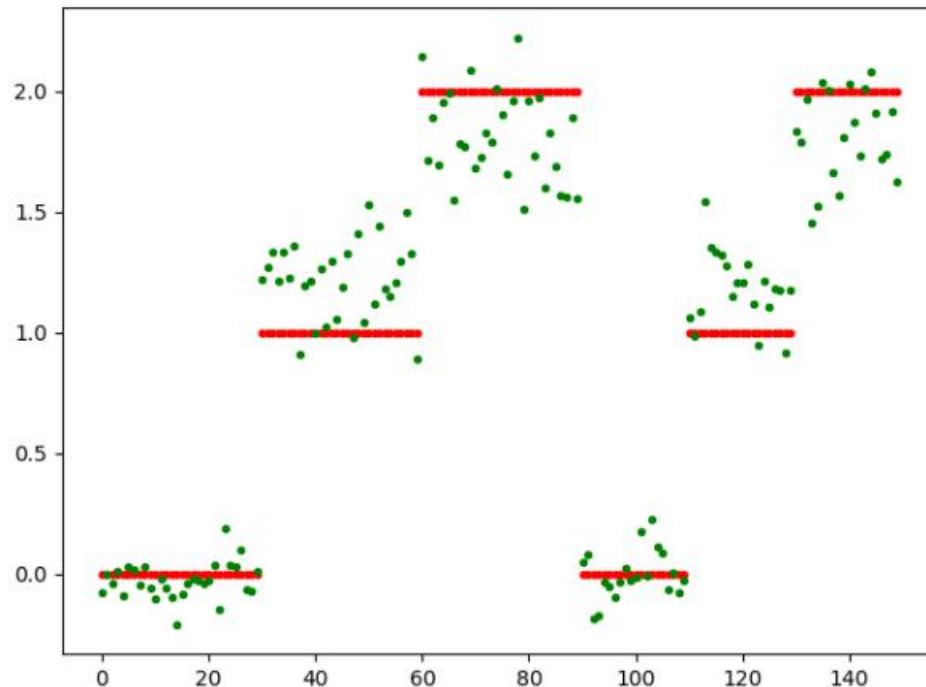
```

The prediction values of theta are: [0.9995957964945576, -0.033807975361895, -0.059610323598668226, 0.32140174738367494, 0.4677042339633092]
Estimate function is: y = 0.9995957964945576 + ( -0.033807975361895 ) * x1 + ( -0.059610323598668226 ) * x2 + 0.32140174738367494 * x3 + 0.4677042339633092 * x4
The loss of the estimation function to the training set is: 0.024938655178665472
The loss of the estimation function to the testing set is 0.021857580973334987
Accuracy rate is: 0.9666666666666667

```

The result of the loss for training data is about 0.025, and the result of the loss for testing data is about 0.022. These two loss results are similar. Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%.

This is the distribution of the tags:



The three on the left are the results of the training set, and the three on the right are the results of the test set. We can see that the distribution of the fitting results is evenly distributed.

(ii) If we exchange training data and testing data, we can get the following results:

```
The prediction values of theta are: [0.99959796494557, 0.07795772771460194, -0.0499987310159661, 0.32616004426867373, 0.389780222150264713]
Estimate function is: y = 0.99959796494557 + ( 0.07795772771460194 ) * x1 + ( -0.0499987310159661 ) * x2 + 0.32616004426867373 * x3 + 0.389780222150264713 * x4
The loss of the estimation function to the training set is: 0.02119064510757971
The loss of the estimation function to the testing set is: 0.029303708816517806
Accuracy rate is: 0.9666666666666667
```

The result of the loss for this time's training data is about 0.021, and the result of the loss for this time's testing data is about 0.029.

Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%, the accuracy rate is same as (i).

(iii) We assign values to the tags, Replace Iris-setosa with 100, replace Iris-versicolor with 101 and replace Iris-virginica with 102

Similar to (i), we changed the values of the tags, we got the following results:

```
The prediction values of theta are: [100.99593944595028, -0.033807975361879544, -0.05961032359867521, 0.3214017473836348, 0.46770423396333244]
Estimate function is: y = 100.99593944595028 + ( -0.033807975361879544 ) * x1 + ( -0.05961032359867521 ) * x2 + 0.3214017473836348 * x3 + 0.46770423396333244 * x4
The loss of the estimation function to the training set is: 0.024946898420099957
The loss of the estimation function to the testing set is: 0.021865824214770218
Accuracy rate is: 0.9666666666666667
```

The result of the loss for training data is about 0.025, and the result of the loss for testing data is about 0.022. These two loss results are similar.

Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%.

(iv) If we exchange training data and testing data, we can get the following results:

```
The prediction values of theta are: [100.99593944595028, 0.0779577277146599, -0.049998731015984635, 0.326160044268633, 0.38978022150262737]
Estimate function is: y = 100.99593944595028 + ( 0.0779577277146599 ) * x1 + ( -0.049998731015984635 ) * x2 + 0.326160044268633 * x3 + 0.38978022150262737 * x4
The loss of the estimation function to the training set is: 0.021198888349015087
The loss of the estimation function to the testing set is 0.029311952057954522
Accuracy rate is: 0.9666666666666667
```

Different from (iii) results, but similar to (ii) results, the result of the loss for this time's training data is about 0.021, and the result of the loss for this time's testing data is about 0.029.

Then I use the method mentioned above to calculate the accuracy, the result is about 96.7%.

(2) We use 75 data items as training data, 75 data item as testing data (just cut from middle).

(i) We assign values to the tags. Replace Iris-setosa with 0, replace Iris-versicolor with 1 and replace Iris-virginica with 2

We used the above mentioned method to update the values of parameters, after we get the values of parameters, we can use testing data to test the accuracy of this linear regression model.

Firstly, I used cost function to calculate the loss of training data and estimate the loss of testing data.

Then I use testing data to calculate the accuracy, the method is to separately calculate distances between the estimated value of y and the three tag values, the tag with the smallest distance is used as the estimated tag of this data item, then compared with the real tag. If the estimated tag is the same as real tag, the prediction is correct, if not, the prediction is wrong.

Results

Using the above method we can get the following results.

```
The prediction values of theta are: [0.3333199321648524, 0.03847883314343642, -0.08865152742261499, 0.1988305609515164, 0.1851593115435026]
Estimate function is: y = 0.3333199321648524 + ( 0.03847883314343642 ) * x1 + ( -0.08865152742261499 ) * x2 + 0.1988305609515164 * x3 + 0.1851593115435026 * x4
The loss of the estimation function to the training set is: 0.0036430323867714756
The loss of the estimation function to the testing set is 0.9224736241918065
Accuracy rate is: 0.0
```

The result of the loss for training data is about 0.004, and the result of the loss for testing data is about 0.922. These two loss results are so different.

Then I use the method mentioned above to calculate the accuracy, the result is about 0.0%.

(ii) If we exchange training data and testing data, we can get the following results:

```
The prediction values of theta are: [1.666599660824263, -0.08904184869642665, -0.10470318591896022, 0.25552143969294, 0.29481769230738947]
Estimate function is: y = 1.666599660824263 + ( -0.08904184869642665 ) * x1 + ( -0.10470318591896022 ) * x2 + 0.25552143969294 * x3 + 0.29481769230738947 * x4
The loss of the estimation function to the training set is: 0.026823087116011393
The loss of the estimation function to the testing set is 0.8978006732583867
Accuracy rate is: 0.0
```

The result of the loss for this time's training data is about 0.027, and the result of the loss for this time's testing data is about 0.898. These two loss results are so different. Then I use the method mentioned above to calculate the accuracy, the result is about 0.0%.

Summary

Now we know the ways to train linear regression model. If we select use the matrix function, we can get the parameters directly, but if we have many features, the

computational cost of this method will be complicated and only applies to linear models.

If we select use gradient descent, we must select learning rate and steps, and we need to iterate continuously to update the values of the parameters.

We can also know that our training set must be comprehensive enough to make accurate predictions.