

# Liner regression and classification

## 1 Data and objective

I used the iris data set. The data set has five attributes; sepal length, sepal\_width, petal\_length, petal\_width and species. I tried both regression and classification. This regression is to predict petal length of iris setosa from sepal length of iris setosa. This classification is binary classification of iris setosa and iris versicolor by sepal length and petal length.

## 2 Method of regression

I will briefly explain the regression method I implemented and applied.

### 2.1 Liner basis function model

When  $y$  is predicted from the observation data  $x$ , linear basis function model can be used;

$$y(x, w) = w^T \phi(x)$$

where  $w = (w_0, \dots, w_{M-1})^T$ ,  $\phi = (\phi_0, \dots, \phi_{M-1})^T$ , and  $\phi_0(x) = 1$ .

The weight  $w$  is decided from training data  $X = (x_1, \dots, x_N)^T$  together with their corresponding values  $t = (t_1, \dots, t_N)^T$ ;

$$w = (\Phi^T \Phi)^{-1} \Phi^T t$$

where

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_{M-1}(x_N) \end{bmatrix}.$$

The basis function  $\phi$  is decided. I tried the three following functions.

- Polynomial (In this case, called linear regression.)

$$\phi_j(x) = x^j$$

- Gaussian basis function

$$\phi_j(x) = \exp \left[ -\frac{(x - \mu_j)^2}{2s^2} \right]$$

- Sigmoid basis function

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

where  $s, \mu$  are decided appropriately, but it may be better to decide  $\mu$  covering the range of  $x$ .

This method to decide the weight is derived from maximum likelihood estimation of data generation probability assuming the data follows gaussian distribution.

## 2.2 Liner basis function model with regularization

Further I tried regularized least squares method when  $\phi$  is polynomial. This method prevents over-fitting. This method decides the weight by the following equation;

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where  $\lambda$  is decided appropriately.

## 2.3 Validation test

10% of the total data is test data, 90% of the rest is training data and 10 % of the rest is validation data. If the dimension of model is too high, model is over-fitting to data. If the dimension of model is too low, model don't predict data well. Validation test find the optimal dimension of model. How to split data is changed serveral times, and the dimension is choosed when the expected value of  $E_{RMS}$  is minimum;

$$E_{RMS} = \sqrt{\frac{\sum_{i=0}^N (y_i - t_i)^2}{N}}$$

where predict data  $\mathbf{y} = (y_1, \dots, y_N)^T$  and test data  $\mathbf{t} = (t_1, \dots, t_N)^T$ .

## 3 Result of regression

When the dimension of model is 5,  $E_{RMS}$  is higher than value I expected. This cause may be that iris data is not suitable for regression. In fact, the associated task of iris data is classification (<http://archive.ics.uci.edu/ml/datasets/Iris>).

In regualrization liner regression,  $E_{RMS}$  of even 10 dimensional model is as low as them of other dimensions.

I decide that the dimension of model is 4 based on the policy of the method I mentioned above.

All algorithms seems to be able to predict test data well.

## 4 Method of classification

I will briefly explain the binary classification method I implemented and applied.

## 4.1 Least squares classification

When the obseravation data  $x$  is classified into classes  $C_1$  and  $C_2$ , the discriminant function is

$$y(x) = W^T x$$

where

$$W = \begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,M-1} \\ w_{1,0} & w_{1,1} & \cdots & w_{M-1,M-1} \end{bmatrix}^T.$$

If  $y_1(x) > y_2(x)$ , the obseravation data  $x$  is classified into  $C_1$ . If  $y_1(x) < y_2(x)$ , the obseravation data  $x$  is classified into  $C_2$ .

The weight  $w$  is decided from training data  $X = (x_1, \dots, x_N)^T$  ( $x = (1, x_1, \dots, x_{M-1})^T$ ) together with their corresponding values  $T = (t_1, \dots, t_N)$  (If  $x \in C_1, t = (1, 0)^T$  and if  $x \in C_2, t = (0, 1)^T$ );

$$W = (X^T X)^{-1} X^T T.$$

This method to decide weight is derived from minimaization of square error between data and the discriminant function.

## 4.2 Perceptron algorithm

When the observation data  $x$  is classified into classes  $C_1$  and  $C_2$ , the perceptron function is

$$y(x) = f(w^T \phi(x))$$

where

$$f(a) = \begin{cases} +1 & a \geq 0(C_1) \\ -1 & a < 0(C_2) \end{cases}$$

If  $y(x) = +1$ , the observation data  $x$  is classified into  $C_1$ . If  $y(x) = -1$ , the observation data  $x$  is classified into  $C_2$ .

The weight  $w$  is decided from training data  $x = (x_1, \dots, x_N)$  together with their corresponding values  $t$  (If  $x \in C_1, t = +1$  and if  $x \in C_2, t = -1$ );

$$w^{(\tau+1)} = w^{(\tau)} + \eta \phi_n t_n.$$

where  $\eta$  is decided appropriately. If all points or most points are classified correctly, iteration is stopped.

This method to decide weight is derived from minimaization of

$$E_p(w) = - \sum_{n \in M} w^T \phi_n t_n.$$

I tried perceptron algorithm in the case of  $\phi(x) = x$ .

## 4.3 Validation test

10% of the total data is test data, 90% of the rest is training data and 10 % of the rest is validation data. Validation test checks if the discriminant function can classify serveral data correctly. How to split data is changed serveral times, and the expected value of the accuracy rate is calculated.

## **5 Result of classification**

The expected value of the accuracy rate in validation data is 1.0. Both algorithms can classify data correctly in test data.