



华南理工大学

South China University of Technology

The Experiment Report of Machine Learning

SCHOOL: SCHOOL OF SOFTWARE ENGINEERING

SUBJECT: SOFTWARE ENGINEERING

Author:

Peng Cheng

Supervisor:

Mingkui Tan

Student ID:

201830020255

Grade:

Undergraduate

2020-10-24

Linear Regression and Stochastic Gradient Descent

Abstract—In this paper, the linear regression problem is transformed into the optimization problem of minimizing the loss function, and parameters in linear regression are solved based on the closed form solution method and the stochastic gradient descent method respectively.

I. INTRODUCTION

This paper aims to attain the solution of parameters in linear regression. Linear regression is the simplest and most used method in regression analysis and prediction, so it is important us to understand and study it further.

II. METHODS AND THEORY

1.1 Linear regression

Linear regression is one of the regression problems. Linear regression assumes that the target value is linearly related to the features, that is, it satisfies a multiple linear equation. By constructing the loss function, the minimum parameters w and b of the loss function are solved. In general, we can express it as follows:

$$f(x; w, b) = w^T x + b$$

2.2 Loss function

The least square loss is used in this paper to measure the performance for regression:

$$L_D(w, b) = 1/n \sum_{i=1}^n (y_i - f(x; w, b))^2$$

And the task is the find the value of W and b when l is minimized, so the optimization formula of the core objective is:

$$w^*, b^* = \operatorname{argmin}_{w, b} L_D(w, b)$$

2.3 Further optimization

In order to simplify the model, we add a regular term to the optimization function with the weights of penalty factor C , and the optimization problem becomes:

$$\min_{w, b} J(w, b) = C \frac{\|w\|^2}{2} + 1/n \sum_{i=1}^n (y_i - f(x; w, b))^2$$

2.4 Solution of parameter

The paper proposed two methods to realize the solution of linear regression parameters: the closed-form solution and Stochastic Gradient Descent (SGD).

In order to simplify our proof, we introduce augmented matrix and augmented vectors still represent w and X .

$$X = (x_1, x_2, \dots, x_n)^T$$

$$w = (b, w_1, \dots, w_n)^T$$

Then the loss function and optimization becomes:

$$L_D(w) = \frac{1}{2} \|y - Xw\|^2$$

$$J_D(w) = \frac{\|w\|^2}{2} + \frac{1}{2} \|y - Xw\|^2$$

The closed-form solution directly derives the loss function and let the derivatives equals 0 to find the Extreme point :

let $(a = y - Xw)$, then

$$\begin{aligned} \frac{\partial L_D(w)}{\partial w} &= \frac{\partial a}{\partial w} \frac{\partial \left(\frac{1}{2} a^T a \right)}{\partial a} \\ &= \frac{1}{2} \frac{\partial a}{\partial w} 2a \\ &= \frac{\partial (y - Xw)}{\partial w} y - X = 0 \end{aligned}$$

Assuming $\|X^T X\| \neq 0$, then we get:

$$w = (X^T X)^{-1} X^T y$$

Another way is SGD. The core content of gradient descent is to update the independent variables constantly (seeking partial derivatives for W and b), so that the objective function is constantly approaching the minimum value.

In this paper, we use $d = -\frac{\partial J_D(w)}{\partial w} = Cw + y - Xw$ as the direction of optimization and update parameters with learning rate η :

$$w' = w - \eta \frac{\partial J_D(w)}{\partial w}$$

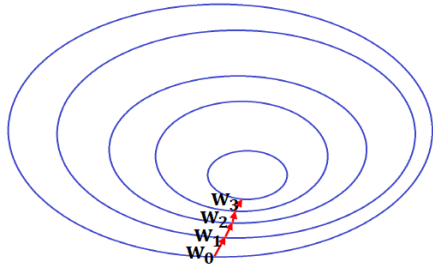


Figure. 1. Sketch map of SGD

Gradient descent uses all samples to calculate the direction, however, information is redundant amongst samples and it presents slow convergence for large data set. So SGD is used in this paper, which randomly select an example in the training set to calculate the direction and we can therefore get to a reasonable solution quickly.

Algorithm: SGD

- 1 Initialize parameter \mathbf{w} and learning rate η
- 2 **While** stopping condition is not achieve
 do
- 3 Randomly select an example i in the
- 4 training set
- 5 $\mathbf{w} = \mathbf{w} - \eta \frac{\partial J_D(\mathbf{w})}{\partial \mathbf{w}}$
- 6 **end**

III. EXPERIMENT

A. Dataset

Linear Regression uses Housing in LIBSVM Data, including 506 samples and each sample has 13 features.

And then it is split into training set and validation set, which comprises 75% and 25% of the raw data, respectively.

B. Implementation

The initialization of the experiment is as follows:

TABLE I

INITIALIZATION OF PARAMETERS

Learning rate	$\eta = 0.005$
Max epoch number	$e = 100$
Penalty factor	$C = 0.5$

And all parameter of \mathbf{w} is initialized randomly.

Using the closed form solution, we get the train loss of 11.2248 and validation loss of 10.6453.

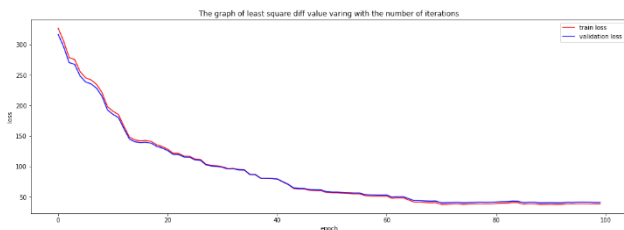


Figure. 2. The graph of least square diff value varying with the number of iterations using SGD

Using SGD, after 100 iterations, the train loss and validation loss drop to 38.1249 and 39.1076, which can be seen from the diagram above:

IV. CONCLUSION

In this experiment, the parameters of linear regression model are solved by two ways. In the closed form solution method, the loss of verification set is 10.6453, while the loss of random gradient descent is 39.1076, which indicates that the fitting degree of the model is high, but meanwhile the loss of random gradient is slightly greater than that of the closed form solution, because the calculated result is not an accurate gradient. For optimization problems, although the loss function obtained by each iteration is not towards the global optimal direction, but the direction of the large whole is towards the global optimal solution, and the final result is often near the global optimal solution.