



The Experiment Report of Machine Learning

SCHOOL: SCHOOL OF SOFTWARE ENGINEERING

SUBJECT: SOFTWARE ENGINEERING

Author:
Changxi Zhu

Supervisor:
Qingyao Wu

Student ID:
201720144894

Grade:
Graduate

December 15, 2017

Logistic Regression, Linear Classification and Stochastic Gradient Descent

Abstract—The experiment consists of two parts, one is logistic regression, the other is linear classification. Both models are updated with four different gradient descent methods and tested in the *a9a* in *LIBSVM Data*. For logistic regression, the accuracy of validation set used NAG, RMSProp, AdaDelta and Adam all can reach 81.598%. For linear classification, the accuracy of validation set used NAG, RMSProp, AdaDelta and Adam all can reach 82.949%.

I. INTRODUCTION

The experiment is divided into two parts, the first is a logistic regression, the second is a linear classification, both are solved by four different gradient descent methods. Through analyzing the different updating process, we can further understand the principle of gradient descent. We will practice on a bigger scale datasets to understand the process of optimization and parameter adjustment. We expect both logistic regression and linear classification can get higher accuracy.

II. METHODS AND THEORY

A. Logistic Regression

For data sets $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$, where $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}; 1)$, $y_i \in \{0, 1\}$, we hope to get a logistic function,

$$f_w(\mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i) \quad (1)$$

where $\mathbf{w} = (w_1; w_2; \dots; w_d; b)$ and

$\sigma(z) = \frac{1}{1 + \exp(-z)}$, to make the $f_w(\mathbf{x}_i)$ got by the model is 1

when it is greater than the threshold and 0 when it is less than the threshold. We hope the results obtained are same with the actual label of \mathbf{x}_i as much as possible.

We use the gradient descent method to solve the model parameters \mathbf{w} . The loss function is

$$-\ln L(\mathbf{w}) = -\frac{1}{K} [\mathbf{y}^T \ln f_w(\mathbf{X}) + (1 - \mathbf{y})^T \ln (1 - f_w(\mathbf{X}))] + \lambda \frac{1}{N} \mathbf{w}^T \mathbf{w} \quad (2)$$

where $\mathbf{X} = (\mathbf{x}_1^T; \mathbf{x}_2^T; \dots; \mathbf{x}_K^T)$, $\mathbf{y} = (y_1; y_2; \dots; y_K)$, λ is a regularization parameters that is artificially adjusted. K is a value between 1 and N , used to implement the stochastic gradient descent.

Our goal is to minimize this loss function, so we need to update the parameters along the negative direction gradient of parameter \mathbf{w} . Find the partial derivative of \mathbf{w} for the loss function, we get the gradient of them

$$G(\mathbf{w}) = \frac{1}{K} (f_w(\mathbf{X}) - \mathbf{y})^T \mathbf{X} + \lambda \frac{2}{K} \mathbf{w} \quad (3)$$

Here we apply four different gradient descent methods to update \mathbf{w} .

NAG:

$$\mathbf{g}^t = G(\mathbf{w}^{t-1} - \gamma \mathbf{v}^{t-1}) \quad (4)$$

$$\mathbf{v}^t = \gamma \mathbf{v}^{t-1} + \eta \mathbf{g}^t \quad (5)$$

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \mathbf{v}^t \quad (6)$$

RMSProp:

$$\mathbf{g}^t = G(\mathbf{w}^{t-1}) \quad (7)$$

$$\mathbf{H}^t = \gamma \mathbf{H}^t + (1 - \gamma) \mathbf{g}^t \odot \mathbf{g}^t \quad (8)$$

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \frac{\eta}{\sqrt{\mathbf{H}^t + \epsilon}} \odot \mathbf{g}^t \quad (9)$$

AdaDelta:

$$\mathbf{g}^t = G(\mathbf{w}^{t-1}) \quad (10)$$

$$\mathbf{H}^t = \gamma \mathbf{H}^t + (1 - \gamma) \mathbf{g}^t \odot \mathbf{g}^t \quad (11)$$

$$\Delta \mathbf{w}^t = -\frac{\sqrt{\Delta_{t-1} + \epsilon}}{\sqrt{\mathbf{H}^t + \epsilon}} \odot \mathbf{g}^t \quad (12)$$

$$\mathbf{w}^t = \mathbf{w}^{t-1} + \Delta \mathbf{w}^t \quad (13)$$

$$\Delta_t = \gamma \Delta_{t-1} + (1 - \gamma) \Delta \mathbf{w}^t \odot \Delta \mathbf{w}^t \quad (14)$$

Adam:

$$\mathbf{g}^t = G(\mathbf{w}^{t-1}) \quad (15)$$

$$\mathbf{m}^t = \beta_1 \mathbf{m}^{t-1} + (1 - \beta_1) \mathbf{g}^t \quad (16)$$

$$\mathbf{H}^t = \gamma \mathbf{H}^t + (1 - \gamma) \mathbf{g}^t \odot \mathbf{g}^t \quad (17)$$

$$\alpha = \eta \frac{\sqrt{1 - \gamma^t}}{1 - \beta^t} \quad (18)$$

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \alpha \frac{\mathbf{m}^t}{\sqrt{\mathbf{H}^t + \epsilon}} \quad (19)$$

where, γ, η, β are super-parameters that are artificially adjusted.

\mathbf{w} is updated by equations (4) to (19) until the loss function converges.

B. Linear classification

For data sets $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, where $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}; 1)$, $y_i \in \{-1, +1\}$. We want to get a linear equation

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i \quad (20)$$

where $\mathbf{w} = (w_1; w_2; \dots; w_d; b)$, to make the $f(\mathbf{x}_i)$ got by the model is +1 when it is greater than the threshold (here 0) and -1 when it is less than the threshold. We hope the results obtained are same with the actual label of \mathbf{x}_i as much as possible.

We also use the gradient descent method to solve the model's parameters. The loss function is

$$\begin{aligned} L(\mathbf{w}, b) &= \sum_{i=1}^K l_{\text{hinge}}(\mathbf{x}_i, y_i) + \lambda \|\mathbf{w}\|_2 \\ &= \frac{1}{K} \left(\sum_{i=1}^K \max(0, 1 - y_i f(\mathbf{x}_i)) + \lambda \|\mathbf{w}\|_2 \right) \end{aligned} \quad (21)$$

where λ is a regularization parameters that is artificially adjusted. K is a value between 1 and N , used to implement the stochastic gradient descent.

Our goal is to minimize this loss function, so we need to update the parameters along the negative direction gradient of parameter \mathbf{w} . Find the partial derivative of \mathbf{w} for the loss function, we get the gradient of them

$$G(w_j) = \frac{\partial L}{\partial w_j} = \frac{1}{K} \sum_{i=1}^K [-\delta(y_i f(\mathbf{x}_i) < 1) y_i x_{ij}] + \frac{1}{K} \lambda \mathbf{w}^T \mathbf{w} w_j \quad (22)$$

Then update \mathbf{w} using the four different gradient descent methods shown in equations (4) to (19) until the loss function converges.

III. EXPERIMENT

A. Dataset

Both experiments use *a9a* of *LIBSVM Data*, including 32561/16281(testing) samples and each sample has 123/123 (testing) features.

B. Implementation

About logistic regression, after reading the dataset, we need to initialize the model, all the parameters used in the experiment are listed in the TABLE I.

TABLE I
PARAMETERS INITIALIZATION

Max Iterations	500
Lambda	0.01
gamma	NAG = 0.9 RMSProp = 0.9 AdaDelta = 0.9 Adam = 0.9
eta	1e-8
beta	Adam = 0.9/0.99

About the \mathbf{w} , initialize them randomly in range(0, 1).

And then, use the model described in II.A to update the \mathbf{w} until the loss function has converged. Fig. 1. And Fig 2 shows the change of loss value with the number of iterations in different gradient descent methods. To my surprise, the curves of the four different gradient descent methods in both two figures almost completely coincide.

About the linear classification, after reading the dataset, we also need to initialize the model, all the parameters used in the experiment are the same as linear regression, shown in Table I.

The \mathbf{w} are also initialized randomly in range(0, 1).

And then, use the model described in II.B to update the \mathbf{w} until the loss function has converged. Fig. 3 and 4. shows the change of loss value with the number of iterations in different gradient descent methods. The same as logistic regression, the curves of the four different gradient descent methods in both two figures almost completely coincide.

Fig. 1-1. Change of loss value with the number of iterations in linear regression in learning rate 0.1 with NGA.

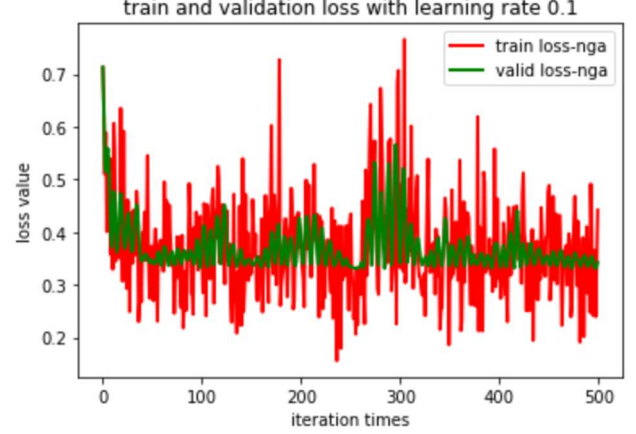


Fig. 1-2. Change of the loss with the number of iterations in the validation set in learning rate 0.01 with NGA .

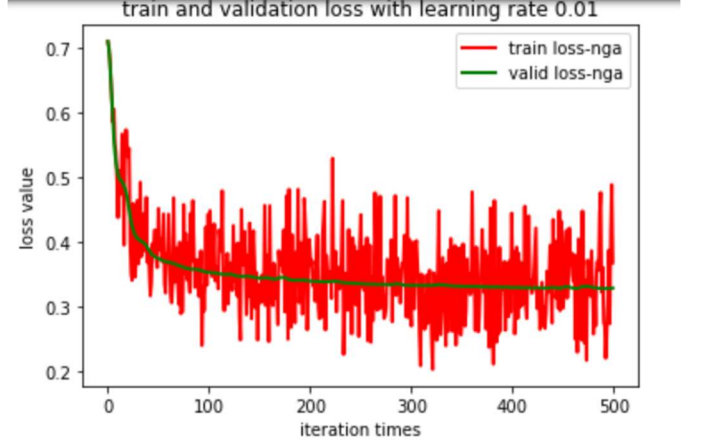


Fig. 1-3. Change of the loss with the number of iterations in the validation set in learning rate 0.001 with NGA .

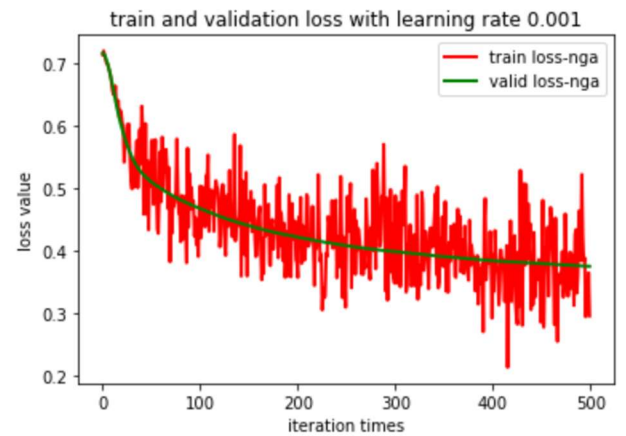


Fig. 2-1. Change of the loss with the number of iterations in the validation set in learning rate 0.1 with rmsprop .

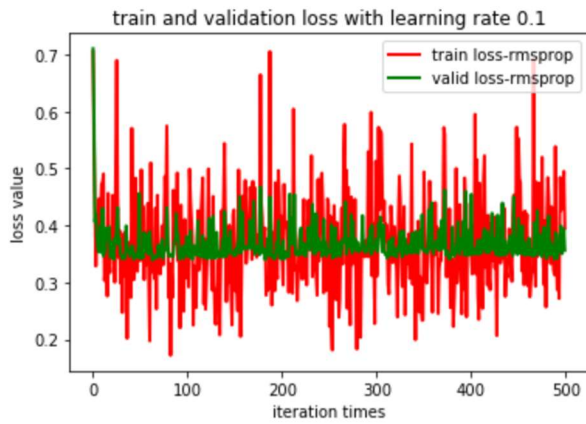


Fig. 2-2. Change of the loss with the number of iterations in the validation set in learning rate 0.01 with rmsprop.

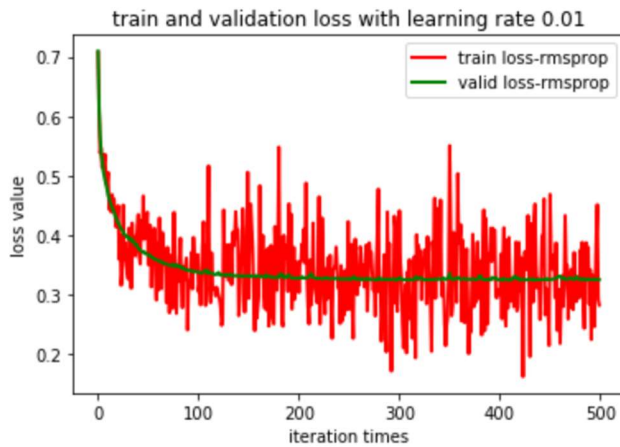


Fig. 2-3. Change of the loss with the number of iterations in the validation set in learning rate 0.001 with rmsprop .

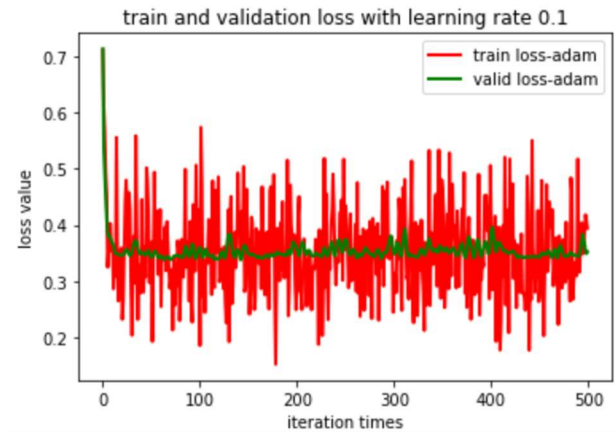
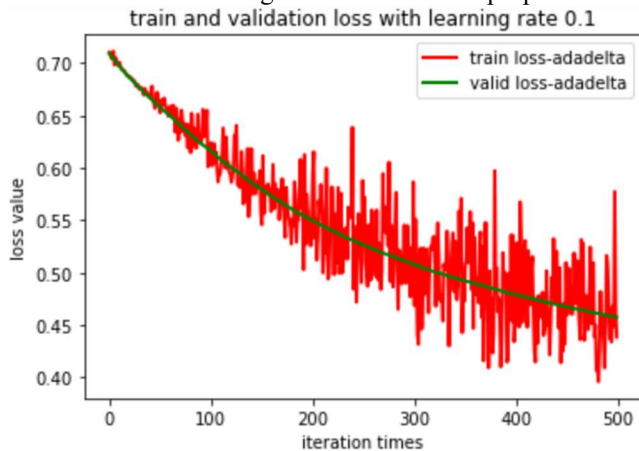


Fig. 3, 4. Change of the loss with the number of iterations in the validation set with svm model
With more pic can see the code

IV. CONCLUSION

The experiment consists of two parts. The first one is logistic regression. The model parameters are updated by four different gradient descent methods. The loss can drop into low level. The second one is a linear classification. Similarly, the model parameters are updated through the same four gradient descent methods. The accuracies of NAG, RMSProp, AdaDelta and Adam on the validation set are all 82.949%.

Through this experiment, further understand the improved version of gradient descent. Through practice in a bigger scale data set, get more experience in adjusting parameters.