



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: MengDan Zheng

Supervisor: Mingkui Tan

Student ID: 201530613795

Grade: Undergraduate

December 15, 2017

Linear Regression, Linear Classification and Gradient Descent

Abstract—Motivation of Experiment is :

① Compare and understand the difference between gradient descent and stochastic gradient descent.

② Compare and understand the differences and relationships between Logistic regression and linear classification.

③ Further understand the principles of SVM and practice on larger data.

I. INTRODUCTION

The experiment is divided into two parts, the first part is to deal with classification problem using logistic regression. The second part is to deal with the classification problem using svm linear classification. Both of these parts use seven gradient optimization algorithms, including SGD, Momentum, NAG, Adagrad, RmsProp, Adadelta and Adam. We will make the comparison of the effects of logistic regression and svm on classification, and the comparison of these optimization algorithms.

II. METHODS AND THEORY

First, about logistic regression, the loss function is:

$$j_i(w) = \sum_{i=1}^n (y_i \log(1/e^{-(w \cdot x_i)}) + (1-y_i) \log(1-1/e^{-(w \cdot x_i)}))$$

The gradient of loss function about w is:

$$\text{grad}_w = (1/e^{-(w \cdot x)} - y) \cdot x.$$

There is one thing to note is the original y is in the set of $\{-1, 1\}$. The formulation above is with the precondition that y is in the set of $\{0, 1\}$. So we should convert the y from -1 to 0 before starting the computation of loss and gradient.

About linear classification, the loss function is:

$$j(w, b) = 1/2 w \cdot T \cdot w + c \cdot \sum_{i=1}^n (\max(0, 1 - y_i \cdot (w \cdot T \cdot x_i + b)))$$

The gradient of loss function about w, b is:

$$\text{If } 1 - y_i \cdot (w \cdot T \cdot x_i + b) > 0 :$$

$$g_{wi} = -y_i \cdot x_i;$$

$$g_{bi} = -y_i;$$

$$\text{if } 1 - y_i \cdot (w \cdot T \cdot x_i + b) < 0 :$$

$$g_{wi} = 0;$$

$$g_{bi} = 0;$$

$$\text{grad}_w = w + c / n \cdot \sum_{i=1}^n (g_{wi});$$

$$\text{grad}_b = c / n \cdot \sum_{i=1}^n (g_{bi})$$

We can let b be a element of w, so we should insert a column with all one to x data. Then the loss function and gradient of linear classification is:

Loss function :

$$j(w) = 1/2 w \cdot T \cdot w + c \cdot \sum_{i=1}^n (\max(0, 1 - y_i \cdot (w \cdot T \cdot x_i)))$$

Gradient of w:

$$\text{If } 1 - y_i \cdot (w \cdot T \cdot x_i) > 0 :$$

$$g_{wi} = -y_i \cdot x_i;$$

$$\text{if } 1 - y_i \cdot (w \cdot T \cdot x_i) < 0 :$$

$$g_{wi} = 0;$$

$$\text{grad}_w = w + c / n \cdot \sum_{i=1}^n (g_{wi});$$

Second, about those gradient optimization algorithms. I learned the fomulation from the blog:

SGD:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J_i(\theta)$$

$$g_t \leftarrow \nabla J_i(\theta_{t-1})$$

$$\theta_t \leftarrow \theta_{t-1} - \eta g_t$$

Momentum:

$$g_t \leftarrow \nabla J(\theta_{t-1})$$

$$v_t \leftarrow \gamma v_{t-1} + \eta g_t$$

$$\theta_t \leftarrow \theta_{t-1} - v_t$$

NAG:

$$g_t \leftarrow \nabla J(\theta_{t-1} - \gamma v_{t-1})$$

$$v_t \leftarrow \gamma v_{t-1} + \eta g_t$$

$$\theta_t \leftarrow \theta_{t-1} - v_t$$

AdaGrad:

$$g_t \leftarrow \nabla J(\theta_{t-1})$$

$$G_t \leftarrow G_t + g_t \odot g_t$$

$$\theta_t \leftarrow \theta_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

RMSProp:

$$g_t \leftarrow \nabla J(\theta_{t-1})$$

$$G_t \leftarrow \gamma G_t + (1 - \gamma) g_t \odot g_t$$

$$\theta_t \leftarrow \theta_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

AdaDelta:

$$g_t \leftarrow \nabla J(\theta_{t-1})$$

$$G_t \leftarrow \gamma G_t + (1 - \gamma) g_t \odot g_t$$

$$\Delta \theta_t \leftarrow - \frac{\sqrt{\Delta_{t-1} + \epsilon}}{\sqrt{G_t + \epsilon}} \odot g_t$$

$$\theta_t \leftarrow \theta_{t-1} + \Delta \theta_t$$

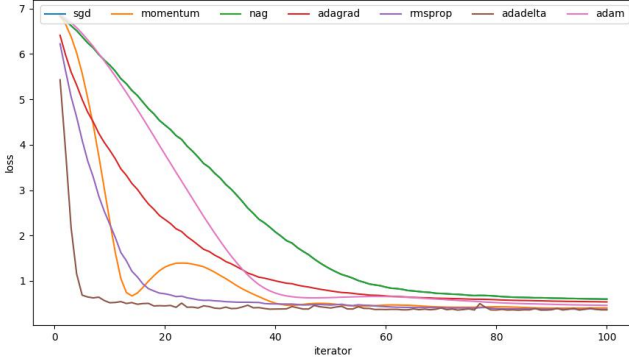
$$\Delta_t \leftarrow \gamma \Delta_{t-1} + (1 - \gamma) \Delta \theta_t \odot \Delta \theta_t$$

Adam:

$$\begin{aligned}
\mathbf{g}_t &\leftarrow \nabla J(\boldsymbol{\theta}_{t-1}) \\
\mathbf{m}_t &\leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \\
G_t &\leftarrow \gamma G_t + (1 - \gamma) \mathbf{g}_t \odot \mathbf{g}_t \\
\alpha &\leftarrow \eta \frac{\sqrt{1 - \gamma^t}}{1 - \beta^t} \\
\boldsymbol{\theta}_t &\leftarrow \boldsymbol{\theta}_{t-1} - \alpha \frac{\mathbf{m}_t}{\sqrt{G_t + \epsilon}}
\end{aligned}$$

III. EXPERIMENT

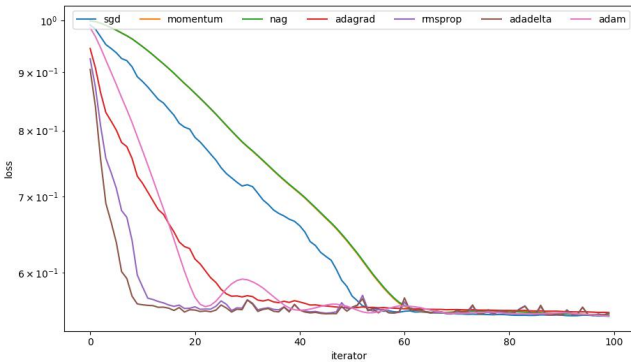
First part is about logistic regression. I used seven gradient optimization algorithms to update the w . The chart of the change of loss is:



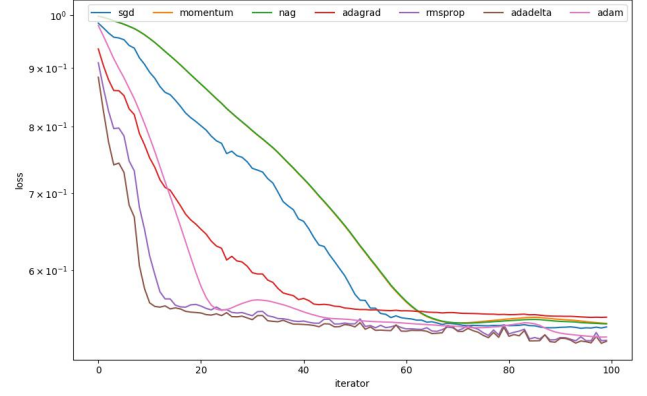
In this section, I set the original learning rate of SGD, Momentum, NAG, RmsProp and Adam are all 0.03. The original rate of adagrad is 0.06 because I found it is too slow to convergence. So in my experiment, adagrad is the slowest algorithm. And adadelta doesn't need original rate.

From the chart, we can get that nag and sgd are very closed. They are the second slowest algorithm to convergence followed by adam, rmsprop, momentum. Adadelta is the fastest algorithm.

Second part is svm linear classifier. In this part I achieved it in two way. First, with the variable w and b , The chart of the change of loss is:



Second, with the only variable w , The chart of the change of loss is:



In this two way, the curves are very similar. I set the original learning rate of Momentum, NAG, RmsProp and Adam are all 0.0005. The original rate of adagrad is 0.01 and rate of sgd is 0.005. The reason is same as the one mentioned above. So the result is so closed with the first part. Adagrad is slower, followed by sgd. NAG and Momentum are very closed. Adadelta, Adam and Rmsprop are all very fast.

As the comparison of logistic regression and svm on classification, I compare them by the accuracy of classification. The classification accuracy of linear regression is about ,and the worst is 0.8304772434125668 . The classification accuracy of svm is about 0.7637737239727289. This indicates logistic regression has a better performance than linear svm on this dataset.

IV. CONCLUSION

From this experiment, It can be concluded that logistic regression has a higher classification accuracy than linear svm. What's more, The order of convergence rate of seven different optimization algorithms from slow to fast are: AdaGrad, SGD, NAG, Momentum, Adam, RMSProp, Adadelta. But the faster the algorithm, the effect is not necessarily better. It should be adjusted according to different circumstances.