



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:
Fangyuan Ding

Supervisor:
Qingyao Wu

Student ID: 201720144979

Grade:
Graduate

December 14, 2017

Logistic Regression, Linear Classification and Stochastic Gradient Descent

Abstract—

Compare and understand the difference between gradient descent and stochastic gradient descent. Compare and understand the differences and relationships between Logistic regression and linear classification. Further understand the principles of SVM and practice on larger data.

I. INTRODUCTION

In the experiment, we use the Logistic Regression and Linear Classification model to solve the classification problem. And use different optimized methods (NAG, RMSProp, AdaDelta and Adam) to train the model in order to Compare and understand the difference between gradient descent and stochastic gradient descent.

II. METHODS AND THEORY

A. Logistic Regression and Stochastic Gradient Descent

1) The loss function of logistic regression

$$\text{Loss} = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

2) Derivatives

$$\frac{\partial \text{Loss}}{\partial w} = -\frac{1}{m} \sum_{i=1}^m \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

B. Linear Classification and Stochastic Gradient Descent

1) The loss function of linear classification

$$\text{Loss} = \frac{\|w\|^2}{2} + \frac{C}{n} \sum_{i=1}^n \max(0, 1 - y_i(x_i w + b))$$

2) Derivatives

$$g_t = \begin{cases} w + C \sum_{i=1}^n -x_i^T y_i & 1 - y_i(x_i w + b) \geq 0 \\ w & 1 - y_i(x_i w + b) < 0 \end{cases}$$

C. Gradient descent optimized methods

1) NAG

$$\begin{aligned} g_t &= \nabla J(\theta_{t-1}) \\ v_t &= \gamma v_{t-1} + \eta g_t \\ \theta_t &= \theta_{t-1} - v_t \end{aligned}$$

2) RMSProp

$$\begin{aligned} g_t &= \nabla J(\theta_{t-1}) \\ G_t &= \gamma G_t + (1 - \gamma) g_t \odot g_t \\ \theta_t &= \theta_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \end{aligned}$$

3) AdaDelta

$$\begin{aligned} g_t &= \nabla J(\theta_{t-1}) \\ G_t &= \gamma G_t + (1 - \gamma) g_t \odot g_t \end{aligned}$$

$$\begin{aligned} \Delta \theta_t &= -\frac{\sqrt{\Delta_{t-1} + \epsilon}}{\sqrt{G_t + \epsilon}} \odot g_t \\ \theta_t &= \theta_{t-1} + \Delta \theta_t \\ \Delta_t &= \gamma \Delta_{t-1} + (1 - \gamma) \Delta \theta_t \odot \Delta \theta_t \end{aligned}$$

4) Adam

$$\begin{aligned} g_t &= \nabla J(\theta_{t-1}) \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ G_t &= \gamma G_t + (1 - \gamma) g_t \odot g_t \\ \alpha &= \eta \frac{\sqrt{1 - \gamma^t}}{1 - \beta^t} \\ \theta_t &= \theta_{t-1} - \alpha \frac{m_t}{\sqrt{G_t - \epsilon}} \end{aligned}$$

III. EXPERIMENT

A. Dataset

In this experiment, we use a9a of LIBSVM Data, including 32561/16281 (testing) samples and each sample has 123/123 (testing) features.

B. Experimental steps

1) Logistic Regression and Stochastic Gradient Descent

1. Load the training set and validation set.
2. Initialize logistic regression model parameters.
3. Select the loss function and calculate its derivation.
4. Calculate gradient G toward loss function from partial samples.
5. Update model parameters using different optimized methods (NAG, RMSProp, AdaDelta and Adam).

6. Select the appropriate threshold, mark the sample whose predict scores greater than the threshold as positive, on the contrary as negative. Predict under validation set and get the different optimized method loss L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ and L_{Adam} .

7. Repeat step 4 to 6 for several times, and drawing graph of L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ and L_{Adam} with the number of iterations.

2) Linear Classification and Stochastic Gradient Descent

1. Load the training set and validation set.
2. Initialize SVM model parameters, you can consider initializing zeros, random numbers or normal distribution.
3. Select the loss function and calculate its derivation, find more detail in PPT.

4. Calculate gradient G toward loss function from partial samples.

5. Update model parameters using different optimized methods (NAG, RMSProp, AdaDelta and Adam).

6. Select the appropriate threshold, mark the sample whose predict scores greater than the threshold as positive, on the contrary as negative. Predict under validation set and get the different optimized method loss L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ and L_{Adam} .

7. Repeat step 4 to 6 for several times, and drawing graph of L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ and L_{Adam} with the number of iterations.

C. initialization model parameters

1) Logistic Regression and Stochastic Gradient Descent

All parameters are set into zero in the linear model.

2) Linear Classification and Stochastic Gradient Descent

All parameters are set into zero in the SVM model.

D. Experimental results and curve

For all gradient descent optimized methods, hyper-parameter were selected as follows:

NAG: $\gamma=0.9$ $\eta=0.5$

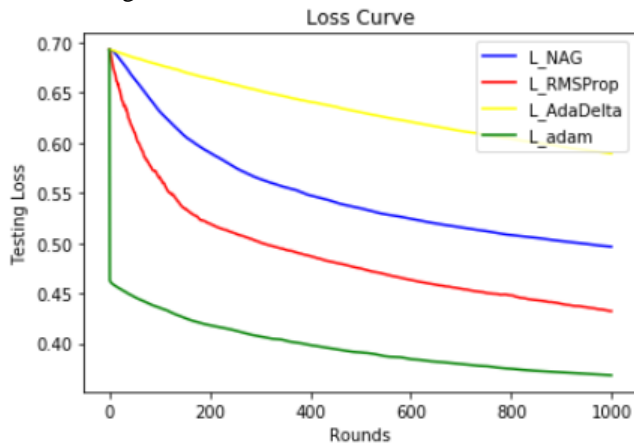
RMSProp: $\gamma=0.9$ $\eta=0.5$

AdaDelta: $\gamma=0.95$ $\Delta t=0$

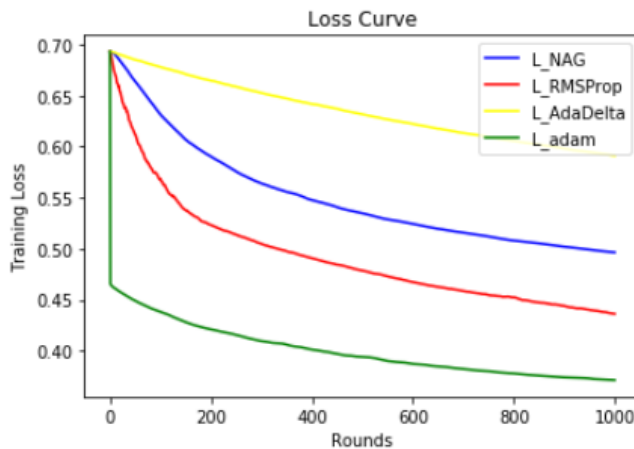
Adam: $\beta t = 0.9$ $\gamma = 0.999$ $\eta = 0.05$

1) Logistic Regression and Stochastic Gradient Descent

Testing Loss:



Training Loss:

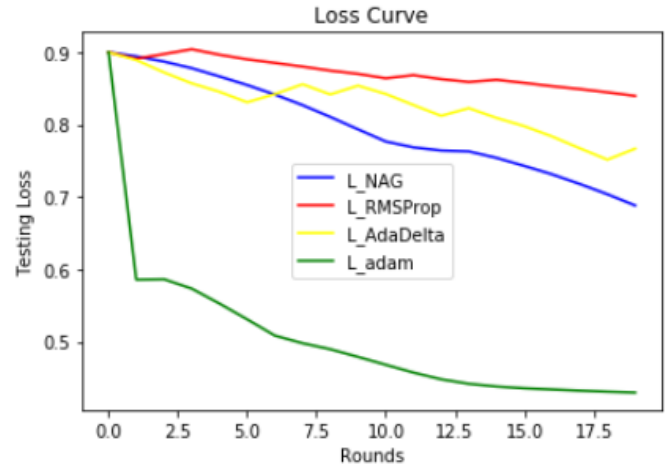


Accuracy:

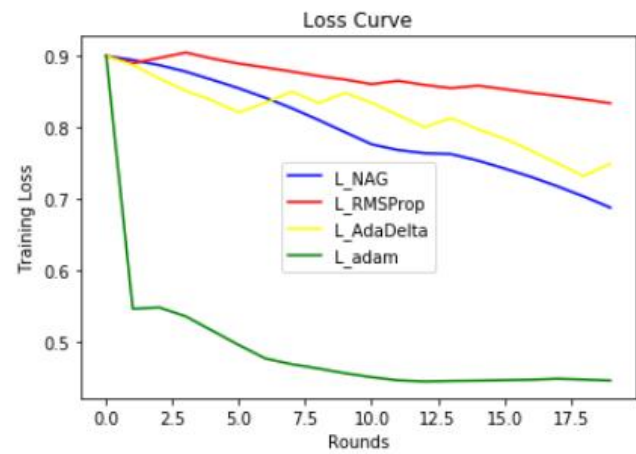
training accuracy_nag= 0.7591904425539756	testing accuracy_nag= 0.7637737239727289
training accuracy_rms= 0.7773102791683302	testing accuracy_rms= 0.7813402125176586
training accuracy_adad= 0.7591904425539756	testing accuracy_adad= 0.7637737239727289
training accuracy_adam= 0.8284450723257886	testing accuracy_adam= 0.8298016092377618

2) Linear Classification Loss:

Testing Loss:



Training Loss:



Accuracy:

training accuracy_nag= 3.071158748195694e-05	testing accuracy_nag= 6.142128861863522e-05
training accuracy_rms= 3.071158748195694e-05	testing accuracy_rms= 6.142128861863522e-05
training accuracy_adad= 3.071158748195694e-05	testing accuracy_adad= 6.142128861863522e-05
training accuracy_adam= 3.071158748195694e-05	testing accuracy_adam= 6.142128861863522e-05

IV. CONCLUSION

Finally, I understand the similarities and the differences between logistic regression and linear classification.

For the similarities, two models solve the classification problem, and both of them are linear models.

For the differences: the logistic regression can represent a probability through mapping features on the sigmoid function. And the linear classification classifies the data by training a hyperplane to split the data.