



The Experiment Report of Machine Learning

SCHOOL: SCHOOL OF SOFTWARE ENGINEERING

SUBJECT: SOFTWARE ENGINEERING

Author:
Shumin Yan

Supervisor:
Mingkui Tan

Student ID:
201530613337

Grade:
Undergraduate

December 14, 2017

Logistic Regression, Linear Classification and Stochastic Gradient Descent

Abstract—In this experiment, we implement logistic regression and SVM using four optimized methods which includes NAG, RMSProp, AdaDelta and Adam. And comparing SVM and logistic regression for solving classification problems.

I. INTRODUCTION

The experiment's purpose is to 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.

In the experiment we need implement logistic regression and SVM using four optimized methods which includes NAG, RMSProp, AdaDelta and Adam.

II. METHODS AND THEORY

As for logistic regression,

$$h_w(x) = g(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

Logistic regression's loss function is

$$J(w) = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log h_w(x_i) + (1 - y_i) \log(1 - h_w(x_i)) \right]$$

For a sample, the gradient function is

$$\begin{aligned} \frac{\partial J(w)}{\partial w} &= -\frac{1}{\partial w} \cdot \partial [y \cdot \log h_w(x) + (1 - y) \log(1 - h_w(x))] \\ &= -y \cdot \frac{1}{h_w(x)} \cdot \frac{\partial h_w(x)}{\partial w} + (1 - y) \cdot \frac{1}{1 - h_w(x)} \cdot \frac{\partial h_w(x)}{\partial w} \\ &= -y \cdot \frac{1}{h_w(x)} \cdot \frac{\partial g(w^T x)}{\partial w} + (1 - y) \cdot \frac{1}{1 - h_w(x)} \cdot \frac{\partial g(w^T x)}{\partial w} \\ &= \left(-\frac{xy}{h_w(x)} + \frac{x(1 - y)}{1 - h_w(x)} \right) \cdot g(w^T x) \cdot [1 - g(w^T x)] \\ &= (h_w(x) - y)x \end{aligned}$$

As for SVM,

Its loss function is

$$L = \max(0, 1 - y_i(w^T x_i + b))$$

Optimization:

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

Derivative:

$$\begin{aligned} g_w(x_i) &= \begin{cases} -y_i x_i & 1 - y_i(w^T x_i + b) \geq 0 \\ 0 & 1 - y_i(w^T x_i + b) < 0 \end{cases} \\ g_b(x_i) &= \begin{cases} -y_i & 1 - y_i(w^T x_i + b) \geq 0 \\ 0 & 1 - y_i(w^T x_i + b) < 0 \end{cases} \\ \frac{\partial f(w, b)}{\partial w} &= w + C \sum_{i=1}^N g_w(x_i) \\ \frac{\partial f(w, b)}{\partial b} &= C \sum_{i=1}^N g_b(x_i) \end{aligned}$$

The four optimized methods includes NAG, RMSProp, AdaDelta, and Adam.

NAG

$$\begin{aligned} \mathbf{g}_t &\leftarrow \nabla J(\theta_{t-1} - \gamma \mathbf{v}_{t-1}) \\ \mathbf{v}_t &\leftarrow \gamma \mathbf{v}_{t-1} + \eta \mathbf{g}_t \\ \theta_t &\leftarrow \theta_{t-1} - \mathbf{v}_t \end{aligned}$$

RMSProp

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

AdaDelta

$$\begin{aligned} \mathbf{g}_t &\leftarrow \nabla J(\theta_{t-1}) \\ G_t &\leftarrow \gamma G_t + (1 - \gamma) \mathbf{g}_t \odot \mathbf{g}_t \\ \Delta \theta_t &\leftarrow -\frac{\sqrt{\Delta_{t-1} + \epsilon}}{\sqrt{G_t + \epsilon}} \odot \mathbf{g}_t \\ \theta_t &\leftarrow \theta_{t-1} + \Delta \theta_t \\ \Delta_t &\leftarrow \gamma \Delta_{t-1} + (1 - \gamma) \Delta \theta_t \odot \Delta \theta_t \end{aligned}$$

Adam

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

III. EXPERIMENT

A.Dataset

Experiment uses a9a of LIBSVMData. The training dataset has 32561 samples and each sample has 123 features. The validation dataset has 16281 samples and each has 123 features.

B.Implementation

Logistic regression

Logistic regression's experimental steps:

1. Load the training set and validation set.
2. Initialize logistic regression model parameters, you can consider initializing zeros.
3. Select the loss function and calculate its derivation.
4. Calculate gradient 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.

Hyper-parameter selection:

SGD:

$$\eta = 0.1$$

NAG:

$$\eta = 0.01$$

$$\gamma = 0.9$$

RMSProp:

$$\eta = 0.01$$

$$\gamma = 0.9$$

$$\epsilon = 1e - 8$$

AdaDelte:

$$\gamma = 0.9$$

$$\epsilon = 1e - 8$$

Adam:

$$\eta = 0.001$$

$$\gamma = 0.999$$

$$\epsilon = 1e - 8$$

$$\beta_1 = 0.9$$

Predicted Results (Best Results):

Accuracy rate:

SGD: 0.832

NAG: 0.812

RMSProp: 0.807

AdaDelte: 0.764

Adam: 0.764

Loss value:

SGD:0.378

NAG:0.391

RMSProp:0.398

AdaDelte:0.523

Adam:0.621

Loss curve:

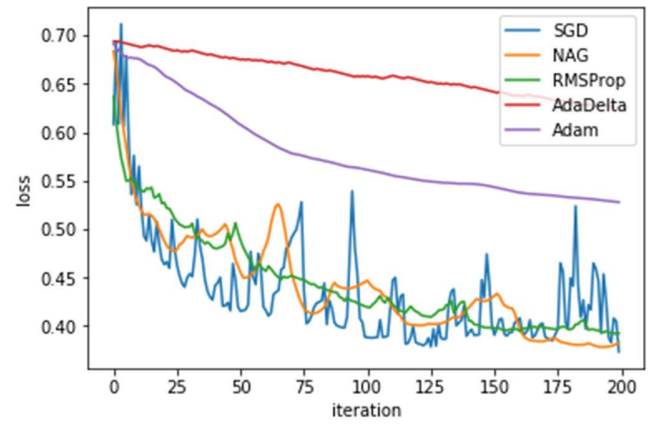


Fig.1 Logistic Regression's loss curve

SVM

SVM's experimental steps:

1. Load the training set and validation set.
2. Initialize SVM model parameters, you can consider initializing zeros.
3. Select the loss function and calculate its derivation.
4. Calculate gradient 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.

Hyper-parameter selection:

SGD:

$$\eta = 0.003$$

$$C = 0.1$$

NAG:

$$C = 0.1$$

$$\eta = 0.01$$

$$\gamma = 0.9$$

RMSProp:

$$\begin{aligned} C &= 0.1 \\ \eta &= 0.001 \\ \gamma &= 0.9 \\ \epsilon &= 1e-8 \end{aligned}$$

AdaDelte:

$$\begin{aligned} C &= 0.1 \\ \gamma &= 0.95 \\ \epsilon &= 1e-8 \end{aligned}$$

Adam:

$$\begin{aligned} C &= 0.1 \\ \eta &= 0.001 \\ \gamma &= 0.999 \\ \epsilon &= 1e-8 \\ \beta_1 &= 0.9 \end{aligned}$$

In conclusion, logistic regression is different from SVM. First, the method of finding the optimal hyperplane is different. Logistic regression finds the hyperplane, to let all points away from it. The hyperplane which SVM looks for, is the most close to the middle line. Second, the SVM can handle the nonlinear case. Third, their loss function is different. The loss of the Logistic regression function is cross entropy loss, and the SVM is the hinge loss.

Predicted Results (Best Results):

Loss value:

SGD: 0.048

NAG: 0.058

RMSProp: 0.051

AdaDelte: 0.052

Adam: 0.047

Loss curve:

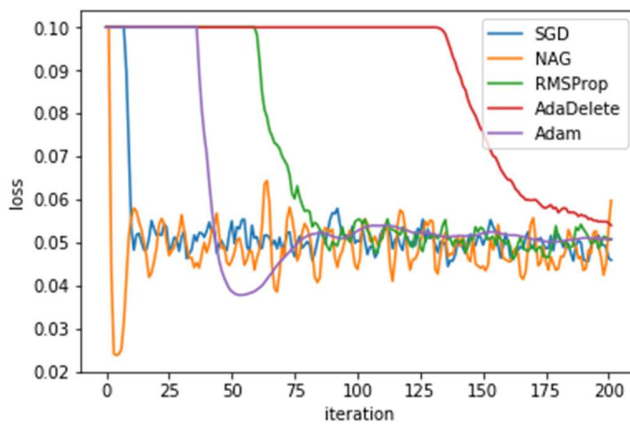


Fig.2 SVM's loss curve

IV. CONCLUSION

After the experiment, I know the differences and relations of the gradient descent, stochastic gradient descent, and batch gradient descent. I implement logic regression and linear classification, and understand their principle and their relationships. In addition, through the optimization methods, I understand what is NAG, RMSProp, AdaDelta and Adam and understand how to use them.