



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:

Fankai Xu

Supervisor:

Qingyao Wu

Student ID: 201721046012

Grade:

graduate

December 15, 2017

# Logistic Regression, Linear Classification and Stochastic Gradient Descent

**Abstract**—Stochastic gradient descent (often shortened to SGD), also known as incremental gradient descent, is a stochastic approximation of the gradient descent optimization and iterative method for minimizing an objective function that is written as a sum of differentiable functions. Compare the differences and relationships between Logistic regression and linear classification.

## I. INTRODUCTION

The purpose of this report is to introduce two kinds of Classifier. In this experiment, we use Logistic Regression Classification and SVM Linear Classification to fit a relatively large data-set. Compare different kind of optimizers like NAG, RMSProp, AdaDelta and Adam, to see characters of each optimizer. Another purpose is to compare Stochastic Gradient Descent and Gradient Descent. In this experiment, the most important thing is to set appropriate parameters (i.e.  $\eta$  the learning rate,  $\gamma$ , epoch,  $\mathbf{W}$  the initial matrix, batch\_size). For the sake of getting a great fit function, we try different learning rate and hyper-parameters for optimizers in the experiment.

## II. METHODS AND THEORY

Logistic regression and SVMs are closely related algorithms, even though this isn't obvious from the usual presentation. In particular, the loss functions are very similar, and linear SVMs and logistic regression can often be substituted for one another without a big difference in performance.

Stochastic Gradient Descent (SGD) is often used in solving optimization problems related to model fitting on huge training data set, searching for the best possible parameters that have minimal cost. Usually, the cost function is continuous and smooth.

When dealing with large data set and performing online (real time) learning, high order gradient descent methods (e.g. Newton's method) is not applicable. Even conjugate gradient descent (CG) is not good because it requires accurate computation of the conjugate descent direction, hence the full data set is necessarily involved in each gradient computation.

On the other hand, using SGD, each time the gradient direction is computed based on only part of the data set (or even a single

data point), so that the computational cost of each gradient step is constant, doesn't scale up with the size of the data set. In addition, it is reported that SGD converges much faster than regular gradient descent (GD).

Optimizer is also an important part of Gradient Descent method.

NAG (Nesterov accelerated gradient) is a way to accelerate the speed of descending. The formula below is how it updates the parameters.

$$\begin{aligned} v_t &= \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1}) \\ \theta &= \theta - v_t \end{aligned} \quad (1)$$

Adadelta is an extension of Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size  $w$ .

RMSprop and Adadelta have both been developed independently around the same time stemming from the need to resolve Adagrad's radically diminishing learning rates. RMSprop in fact is identical to the first update vector of Adadelta that we derived above:

$$\begin{aligned} E[g^2]_t &= 0.9 E[g^2]_{t-1} + 0.1 g_t^2 \\ \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \end{aligned} \quad (2)$$

RMSprop as well divides the learning rate by an exponentially decaying average of squared gradients. Hinton suggests  $\gamma$  to be set to 0.9, while a good default value for the learning rate  $\eta$  is 0.001.

Adaptive Moment Estimation (Adam) is another method that computes adaptive learning rates for each parameter. In addition to storing an exponentially decaying average of past squared gradients  $v_t$  like Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients  $m_t$ , similar to momentum:

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \end{aligned} \quad (3)$$

$m_t$  and  $v_t$  are estimates of the first moment (the mean) and the second moment (the uncentered variance) of the gradients respectively, hence the name of the method. As  $m_t$  and  $v_t$  are initialized

ialized as vectors of 0's, the authors of Adam observe that they are biased towards zero, especially during the initial time steps, and especially when the decay rates are small (i.e.  $\beta_1$  and  $\beta_2$  are close to 1). They counteract these biases by computing bias-corrected first and second moment estimates:

$$\begin{aligned}\hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t}\end{aligned}\quad (4)$$

Formula below is Adam's update rule:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \quad (5)$$

### III. EXPERIMENT

#### A. Dataset

Experiment uses a9a of LIBSVM Data, including 32561/16281(testing) samples and each sample has 123/123 (testing) features. Please download the training set and validation set.

#### B. Implementation

In Logistic Regression experiment, first we loaded a9a dataset for training and a9at dataset for validation, and then changed format of the data and reshaped it. We set epoch and learning rate with proper value, then initialized logistic regression model parameters. 10000 or 5000 data were chosen as inputs of stochastic gradient descent. Then computed the loss values. We choosed four kind of optimizers for optimizing the compute of loss value and model parameters (NAG, RMSProp, AdaDelta, Adam). Finally, we plot a graph (Fig.1 and Fig.2) to show how loss value change with iteration of each epoch. It's all about how the experiment proceed.

It's strange that when we choosed AdaDelta as optimizer, the loss value would become NaN after epoches of training. We still do not figure it out.

In SVM Linearer Classification experiment, we did almost the same works, the only different is that we choosed different loss function. In Logistic Regression experiment, the loss function is cross-entropy loss function, but in SVM experiment, we choosed hing-loss function.

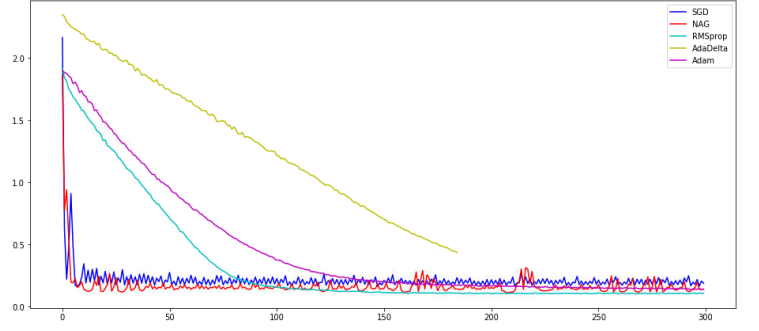


Fig.1. Loss values tendency graph for SGD, NAG, RMSProp, AdaDelta, Adam in SVM experiment.

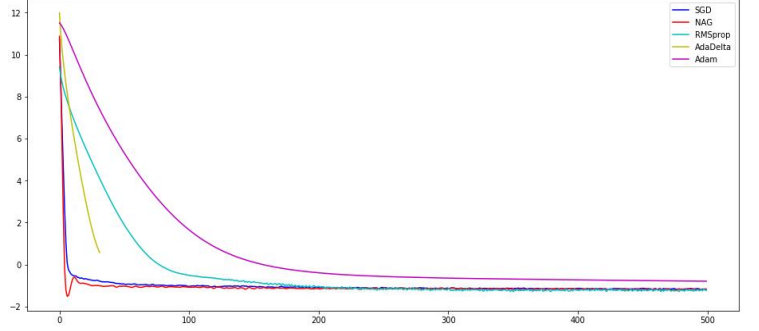


Fig.2. Loss values tendency graph for SGD, NAG, RMSProp, AdaDelta, Adam in Logistic Classification experiment.

### IV. CONCLUSION

From these two experiments, we can see that SVM Linear Regression and Logistic Classification can do the same classification job for not that large dataset. Stochastic Gradient Descent is faster than Gradient Descent in training, because it uses less data to train in each epoch. For each optimizer, we can see from the figures, they have different loss change rates, but all did a good job.