



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:  
Qiaoyuan Wen

Supervisor:  
Qingyao Wu

Student ID:  
201530612989

Grade:  
Undergraduate

December 9, 2017

# Logistic Regression, Linear Classification and Stochastic Gradient Descent

**Abstract—**

## I. INTRODUCTION

In statistics, logistic regression is a regression model where the dependent variable (DV) is categorical. This article covers the case of a binary dependent variable—that is, where the output can take only two values, "0" and "1", which represent outcomes such as pass/fail, win/lose, alive/dead or healthy/sick. Cases where the dependent variable has more than two outcome categories may be analysed in multinomial logistic regression, or, if the multiple categories are ordered, in ordinal logistic regression. In the terminology of economics, logistic regression is an example of a qualitative response/discrete choice model.

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. In other words, SGD tries to find minima or maxima by iteration.

## II. METHODS AND THEORY

### NAG

Further proposals include the NAG method, Stochastic gradient descent with momentum remembers the update  $\Delta w$  at each iteration, and determines the next update as a linear combination of the gradient and the previous update:

$$\Delta w := \alpha \Delta w - \eta \nabla_{\mathbf{w}} Q(\mathbf{w} - \gamma \Delta w)$$

$$\mathbf{w} := \mathbf{w} + \Delta w$$

that leads to:

$$\mathbf{w} := \mathbf{w} - \eta \nabla_{\mathbf{w}} Q(\mathbf{w}) + \alpha \Delta w$$

where the parameter  $\mathbf{w}$  which minimizes  $Q(\mathbf{w})$  is to be estimated, and  $\eta$  is a step size (sometimes called the learning rate in machine learning).

The name momentum stems from an analogy to momentum in physics: the weight vector  $\mathbf{w}$ , thought of as a particle traveling through parameter space, incurs acceleration from the gradient of the loss ("force"). Unlike in classical stochastic gradient descent, it tends to keep traveling in the same direction, preventing oscillations.

### RMSProp

RMSProp (for Root Mean Square Propagation) is also a method in which the learning rate is adapted for each of the parameters. The idea is to divide the learning rate for a weight by a running average of the magnitudes of recent gradients for that

weight.[19] So, first the running average is calculated in terms of means square,

$$v(\mathbf{w}, t) := \gamma v(\mathbf{w}, t-1) + (1-\gamma) (\nabla_{\mathbf{w}} Q(\mathbf{w}))^2$$

Where  $\gamma$  is the forgetting factor. And the parameters are updated as,

$\mathbf{w} := \mathbf{w} - \frac{\eta}{\sqrt{v(\mathbf{w}, t)}} \nabla_{\mathbf{w}} Q(\mathbf{w})$   
RMSProp has shown excellent adaptation of learning rate in different applications. RMSProp can be seen as a generalization of Rprop and is capable to work with mini-batches as well opposed to only full-batches.

### AdaDelta

$$\begin{aligned} & \nabla_{\mathbf{w}} J(\theta_{t-1}) \parallel \mathbf{G}_t \text{ \& } \\ & \text{gets } \gamma \mathbf{G}_t + (1-\gamma) \nabla_{\mathbf{w}} J(\theta_{t-1}) \parallel \mathbf{G}_t \text{ \& } \\ & \text{gets } - \frac{\nabla_{\mathbf{w}} J(\theta_{t-1}) \parallel \mathbf{G}_t}{\sqrt{\mathbf{G}_t \parallel \mathbf{G}_t + \epsilon}} \parallel \nabla_{\mathbf{w}} J(\theta_{t-1}) \text{ \& } \\ & \text{gets } \gamma \Delta_{t-1} + (1-\gamma) \Delta_{t-1} \parallel \nabla_{\mathbf{w}} J(\theta_{t-1}) \text{ \& } \\ & \text{gets } \gamma \Delta_{t-1} + (1-\gamma) \Delta_{t-1} \parallel \nabla_{\mathbf{w}} J(\theta_{t-1}) \end{aligned}$$

### Adam

Adam (short for Adaptive Moment Estimation) is an update to the RMSProp optimizer. In this optimization algorithm, running averages of both the gradients and the second moments of the gradients are used. Given parameters  $\mathbf{w}^{(t)}$  and a loss function  $L^{(t)}$ , where  $t$  indexes the current training iteration (indexed at 1), Adam's parameter update is given by:

$$\mathbf{m}_{\mathbf{w}}^{(t+1)} \leftarrow \beta_1 \mathbf{m}_{\mathbf{w}}^{(t)} + (1-\beta_1) \nabla_{\mathbf{w}} L^{(t)}$$

$$\mathbf{v}_{\mathbf{w}}^{(t+1)} \leftarrow \beta_2 \mathbf{v}_{\mathbf{w}}^{(t)} + (1-\beta_2) (\nabla_{\mathbf{w}} L^{(t)})^2$$

$$\hat{\mathbf{m}}_{\mathbf{w}} = \frac{\mathbf{m}_{\mathbf{w}}^{(t+1)}}{1-\beta_1^{(t)}}$$

$$\hat{\mathbf{v}}_{\mathbf{w}} = \frac{\mathbf{v}_{\mathbf{w}}^{(t+1)}}{1-\beta_2^{(t)}}$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \frac{\hat{\mathbf{m}}_{\mathbf{w}}}{\sqrt{\hat{\mathbf{v}}_{\mathbf{w}} + \epsilon}}$$

where  $\epsilon$  is a small number used to prevent division by 0, and  $\beta_1$  and  $\beta_2$  are the forgetting factors for gradients and second moments of gradients, respectively.

## III. EXPERIMENT

### Motivation of Experiment

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.

### 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.

### Environment for Experiment

python3, at least including following python package: sklearn, numpy, jupyter, matplotlib

It is recommended to install anaconda3 directly, which has built-in python package above.

#### Experiment Step

The experimental code and drawing are completed on jupyter.

#### Logistic Regression and Stochastic Gradient Descent

Load the training set and validation set.

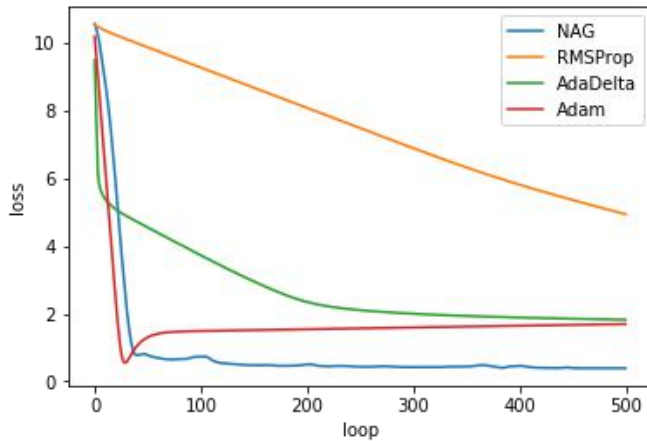
Initialize logistic regression model parameters, you can consider initializing zeros, random numbers or normal distribution.

Select the loss function and calculate its derivation, find more detail in PPT.

Calculate gradient toward loss function from partial samples.

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

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



#### Linear Classification and Stochastic Gradient Descent

Load the training set and validation set.

Initialize SVM model parameters, you can consider initializing zeros, random numbers or normal distribution.

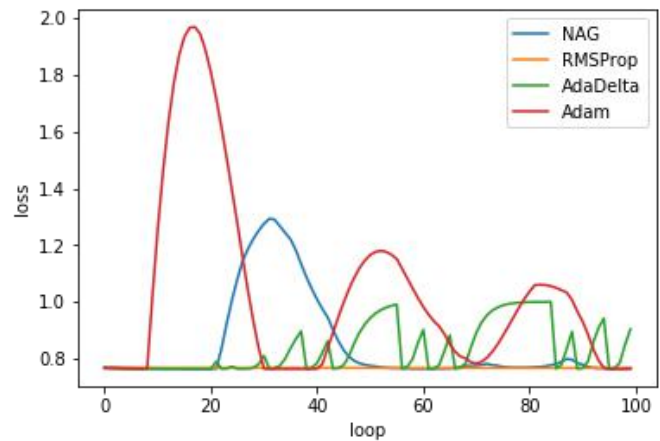
Select the loss function and calculate its derivation, find more detail in PPT.

Calculate gradient toward loss function from partial samples.

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

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

Result:



#### IV. CONCLUSION

Through this experiment, I have a deeper understanding of Logistic Regression, Linear Classification and Stochastic Gradient Descent. For different tasks, you should try different optimization algorithms. Different optimization algorithms have different learning speed, and the final convergence is not necessarily the same. In some cases SGD may also be the best choice.