



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:*  
Shen Fu

*Supervisor:*  
Qingyao Wu

*Student ID:*  
201720144955

*Grade:*  
Graduate

December 15, 2017

# Logistic Regression, Linear Classification and Stochastic Gradient Descent

**Abstract**—Gradient descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks. At the same time, every state-of-the-art Deep Learning library contains implementations of various algorithms to optimize gradient descent. This report aims to introduce the basics of SGD, NAG, Adadelta, RMSprop and Adam. The experiments part will show the result of logistic regression and linear classification using different optimization algorithms.

## I. INTRODUCTION

**G**RAIENT descent optimization algorithms, while increasingly popular, are often used as black-box optimizers, as practical explanations of their strengths and weaknesses are hard to come by. At the same time, every state-of-the-art Deep Learning library contains implementations of various algorithms to optimize gradient descent. These algorithms, however, are often used as black-box optimizers, as practical explanations of their strengths and weaknesses are hard to come by [1].

This experiment aims provide the intuitions towards the behavior of different algorithms for optimizing gradient descent. The first part comes to the introduction of the most common optimization algorithms, content includes the motivation to resolve challenges and how those algorithms lead to the derivation of their update rules. To demonstrate the strengths and weaknesses of different optimization algorithm, a variety of experimental results and analysis that compares those five optimization algorithms with each other will be shown in experiments part. Finally, we will give a conclusion to summarize this experiment.

The experiments are run under the environment of python3. Python package including sklearn, numpy, jupyter, matplotlib are required. Therefore, the anaconda3 are used directly that has built-in python package above. The experimental code and drawing are completed on jupyter.

The rest parts of report are structured as follows: method and theory used in two experiments are presented in Section II. The overall implementation and results of experiments are shown in Section III and Section IV concludes the reports.

## II. METHODS AND THEORY

### A. Stochastic Gradient Descent

Stochastic gradient descent (SGD) in contrast performs a parameter update for each training example  $x^i$  and  $y^i$  label :

$$\theta = \theta - \nabla_{\theta} J(\theta; x^i; y^i) \quad (1)$$

Batch gradient descent performs redundant computations for large datasets, as it recomputed gradients for similar examples before each parameter update. SGD does away with this

redundancy by performing one update at a time. It is therefore usually much faster and can also be used to learn online.

SGD performs frequent updates with a high variance that cause the objective function to fluctuate heavily as in Fig. 1.

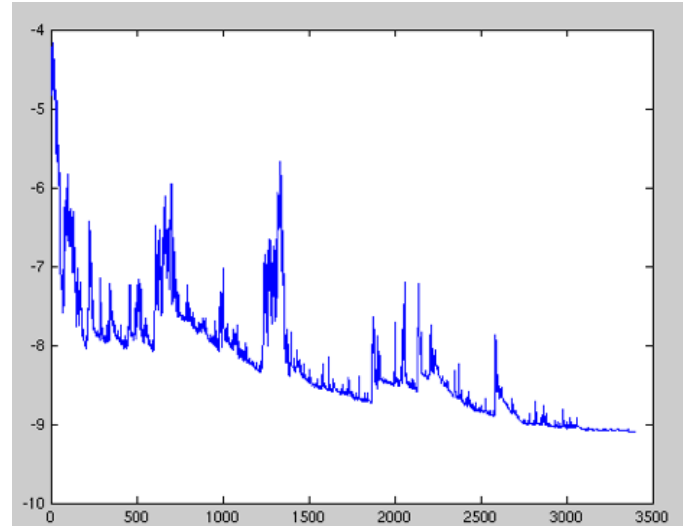


Fig. 1. SGD fluctuation

While batch gradient descent converges to the minimum of the basin the parameters are placed in, SGD's fluctuation, on the one hand, enables it to jump to new and potentially better local minima. On the other hand, this ultimately complicates convergence to the exact minimum, as SGD will keep overshooting. However, it has been shown that when we slowly decrease the learning rate, SGD shows the same convergence behavior as batch gradient descent, almost certainly converging to a local or the global minimum for non-convex and convex optimization respectively. Its code fragment simply adds a loop over the training examples and evaluates the gradient w.r.t. each example.

### B. Momentum

SGD has trouble navigating ravines, i.e. areas where the surface curves much more steeply in one dimension than in another [2], which are common around local optima. In these scenarios, SGD oscillates across the slopes of the ravine while only making hesitant progress along the bottom towards the local optimum as in Fig. 2. Momentum [3] is a method that helps accelerate SGD in the relevant direction and dampens oscillations as can be seen in Fig. 3. It does  $\gamma$  this by adding a fraction of the update vector of the past time step to the current update vector:



Fig. 2. SGD without momentum

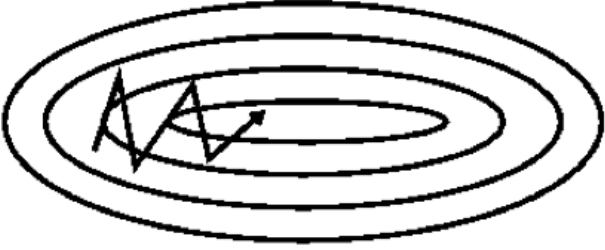


Fig. 3. SGD with momentum

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta) \quad (2)$$

$$\theta = \theta - v_t \quad (3)$$

the momentum term  $\gamma$  is usually set to 0.9 or a similar value.

Essentially, when using momentum, we push a ball down a hill. The ball accumulates momentum as it rolls downhill, becoming faster and faster on the way (until it reaches its terminal velocity if there is air resistance, i.e.  $\gamma < 1$ ). The same thing happens to our parameter updates: The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.

### C. Nesterov Accelerated Gradient

However, a ball that rolls down a hill, blindly following the slope, is highly unsatisfactory. We'd like to have a smarter ball, a ball that has a notion of where it is going so that it knows to slow down before the hill slopes up again.

Nesterov accelerated gradient (NAG) [4] is a way to give our momentum term this kind of prescience. We know that we will use our momentum term  $\gamma v_{t-1}$  to move the parameters  $\theta$ . Computing  $\theta - \gamma v_{t-1}$  thus gives us an approximation of the next position of the parameters (the gradient is missing for the full update), a rough idea where our parameters are going to be. We can now effectively look ahead by calculating the gradient not w.r.t. to our current parameters  $\theta$  but w.r.t. the approximate future position of our parameters:

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1}) \quad (4)$$

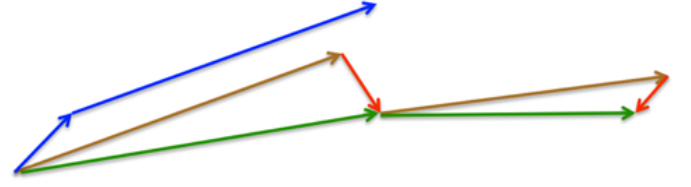


Fig. 4. Nesterov update

$$\theta = \theta - v_t \quad (5)$$

Again, we set the momentum term  $\theta$  to a value of around 0.9. While Momentum first computes the current gradient (small blue vector in Fig. 4) and then takes a big jump in the direction of the updated accumulated gradient (big blue vector), NAG first makes a big jump in the direction of the previous accumulated gradient (brown vector), measures the gradient and then makes a correction (red vector), which results in the complete NAG update (green vector). This anticipatory update prevents us from going too fast and results in increased responsiveness, which has significantly increased the performance of RNNs on a number of tasks [5].

### D. Adadelta

Adadelta [6] 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  $w$ .

Instead of inefficiently storing  $w$  previous squared gradients, the sum of gradients is recursively defined as a decaying average of all past squared gradients. The running average  $E[g^2]_t$  at time step  $t$  then depends (as a fraction  $\gamma$  similarly to the Momentum term) only on the previous average and the current gradient:

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2 \quad (6)$$

We set  $\gamma$  to a similar value as the momentum term, around 0.9. For clarity, we now rewrite our vanilla SGD update in terms of the parameter update vector  $\Delta\theta_t$ :

$$\Delta\theta_t = -\eta \cdot g_{t,i} \quad (7)$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t \quad (8)$$

The parameter update vector of Adagrad that we derived previously thus takes the form:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \varepsilon}} \odot g_t \quad (9)$$

We now simply replace the diagonal matrix  $G_t$  with the decaying average over past squared gradients  $E[g^2]_t$ :

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \varepsilon}} g_t \quad (10)$$



'lambda': [0.01, 0.1], 'eta': [0.01, 0.05], 'Adam\_beta1': [0.9, 0.95], 'Adam\_beta2' : [0.99, 0.999], 'threshold': [0.5, 0.6], 'GD' : 'lambda': [0.01, 0.1, 0.5], 'eta': [0.1, 0.2, 0.3, 0.4, 0.5], 'threshold': [0.4, 0.5, 0.6]

2) *predict result*: The performance of best estimators with five different optimize algorithms of logistic regression and stochastic gradient descent are shown in Table I. And the performance of best estimators with five different optimize algorithms of linear classification and stochastic gradient descent are shown in Table II.

Loss curve of experiment 1 and 2 are shown in Fig.5 and Fig.16, respectively.

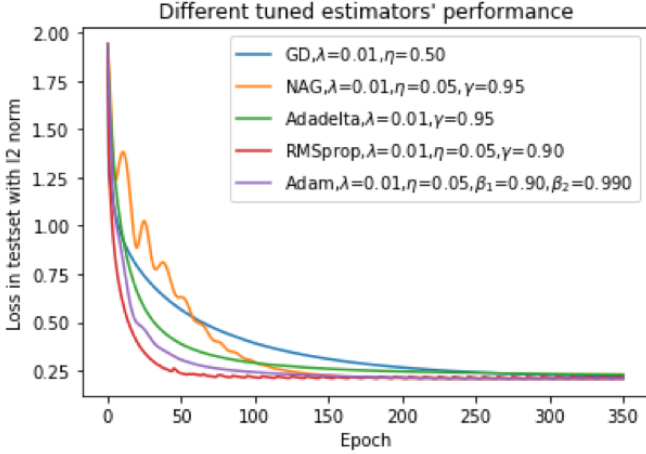


Fig. 5. Loss curve of logic regression

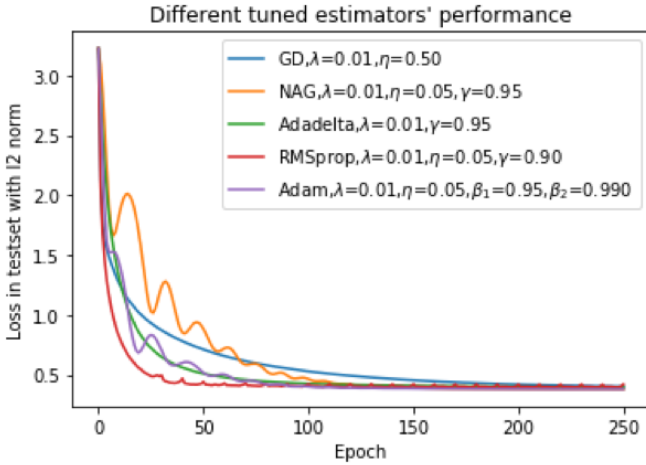


Fig. 6. Loss curve of logic classification

Metrics evolution process details of experiment 1 and 2 are shown in Figs.7-11 and Figs.12-16, respectively.

3) *Analysis*: According to results obtained by two experiments, it can be shown that the performance of different optimization algorithms are ranked as RMSprop, Adam, Adadelta, GD, NAG with exhaustive grid search tuning hyper-parameters. If we can not tune the hyper-parameters due to the resource limits in practice, we can propose RMSprop or Adam.

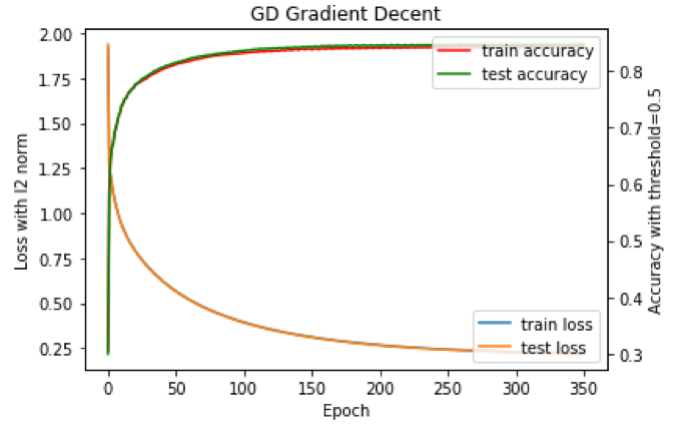


Fig. 7. Metrics evolution process details

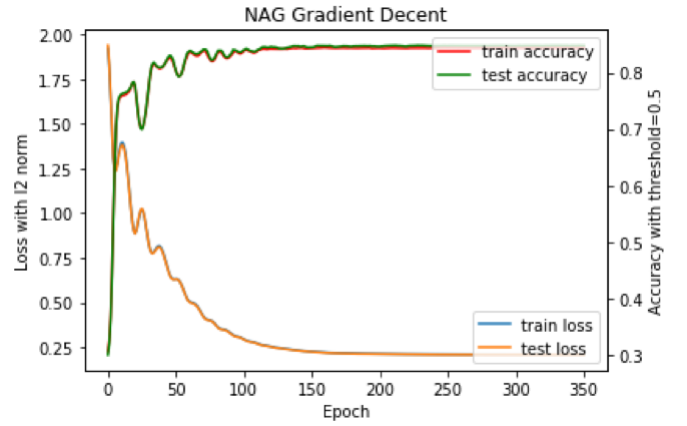


Fig. 8. Metrics evolution process details

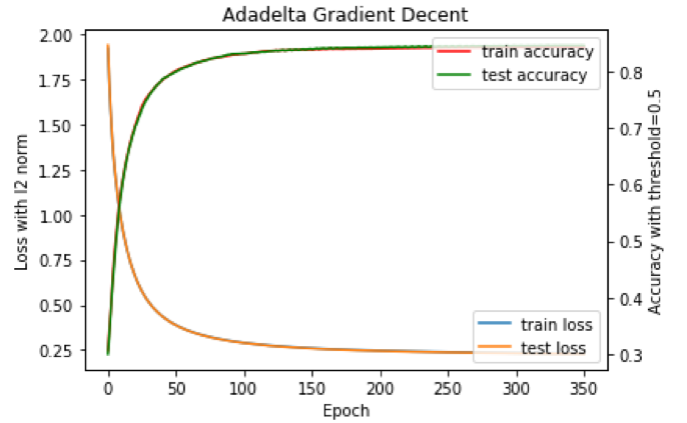


Fig. 9. Metrics evolution process details

Both of the algorithms are not so sensitive to hyper-parameters especially learning rate.

#### IV. CONCLUSION

This report shows two experiments about logic regression and logic classification, respectively. This experiment aims to understand the different between logistic regression and linear

TABLE I

BEST ESTIMATORS WITH FIVE DIFFERENT OPTIMIZE ALGORITHMS OF LOGISTIC REGRESSION AND STOCHASTIC GRADIENT DESCENT

Algorithm	Test accuracy	Train accuracy	Hyper-parameters
GD	0.841927	0.842265	'eta': 0.5, 'lamda': 0.01, 'threshold': 0.5
Adadelata	0.842357	0.842189	'gamma': 0.95, 'lamda': 0.01, 'threshold': 0.5
Adam	0.844139	0.844477	'Adam_beta1': 0.9, 'Adam_beta2': 0.99, 'eta': 0.05, 'lamda': 0.01, 'threshold': 0.5
NAG	0.843954	0.844369	'eta': 0.05, 'gamma': 0.95, 'lamda': 0.01, 'threshold': 0.5
RMSprop	0.844354	0.844139	'eta': 0.05, 'gamma': 0.9, 'lamda': 0.01, 'threshold': 0.6

TABLE II

BEST ESTIMATORS WITH FIVE DIFFERENT OPTIMIZE ALGORITHMS OF LINEAR CLASSIFICATION AND STOCHASTIC GRADIENT DESCEN

Algorithm	Test accuracy	Train accuracy	Hyper-parameters
GD	0.8313319615	0.8317465648	'eta': 0.5, 'lamda': 0.01, 'threshold': 0.4
Adadelata	0.8255888946	0.8259728123	'gamma': 0.95, 'lamda': 0.01, 'threshold': 0.5
Adam	0.8307177298	0.8315162541	'Adam_beta1': 0.95, 'Adam_beta2': 0.99, 'eta': 0.05, 'lamda': 0.01, 'threshold': 0.5
NAG	0.8268480698	0.827800148	'eta': 0.05, 'gamma': 0.95, 'lamda': 0.01, 'threshold': 0.5
RMSprop	0.8290285925	0.8297654828	'eta': 0.05, 'gamma': 0.9, 'lamda': 0.01, 'threshold': 0.5

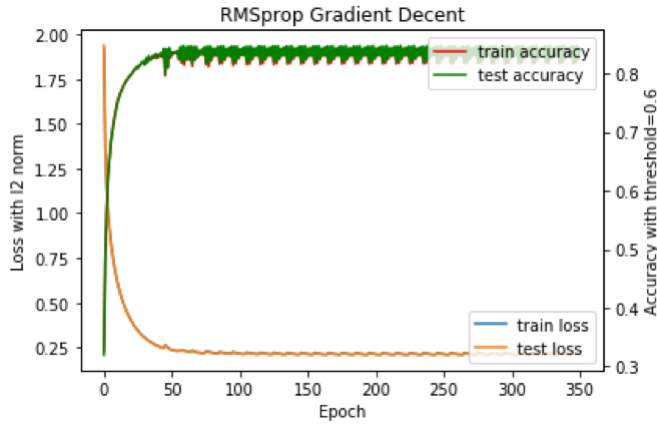


Fig. 10. Metrics evolution process details

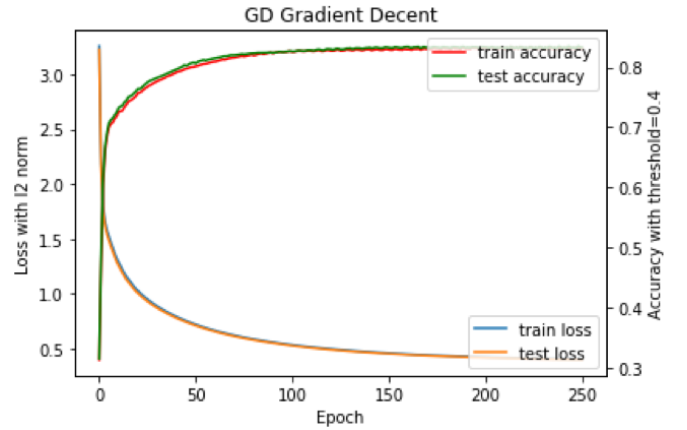


Fig. 12. Metrics evolution process details



Fig. 11. Metrics evolution process details

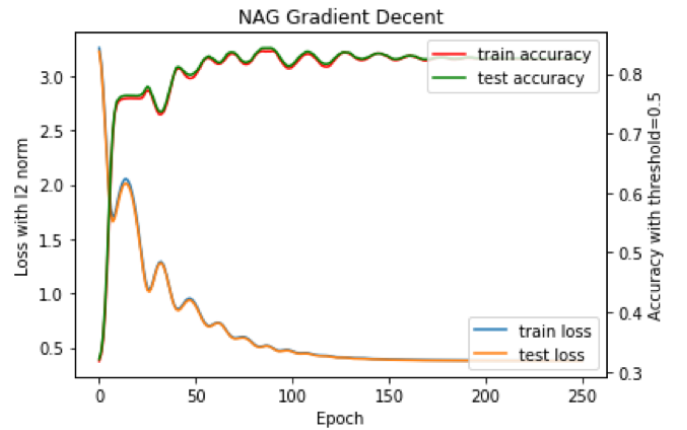


Fig. 13. Metrics evolution process details

classification and learn 5 optimization algorithms: SGD, NAG, Adadelata, RMSprop and Adam. The experiments theroy, implementation are introduced and the detailed results are shown as tables and figures. We can find that logistic regression seems very similar to classification, but the logistic regression pays more attention on its regression operation where the regression object is sigmoid function. Logistic regression and SVM are

both very widely used linear classifier. Furthermore, SGD, NAG, Adadelata, RMSprop and Adam are all widely used methods to optimize gradient decent, while SGD is the most simple method.

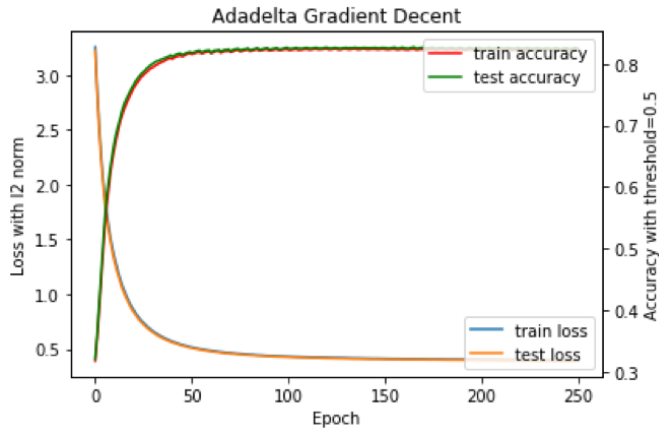


Fig. 14. Metrics evolution process details

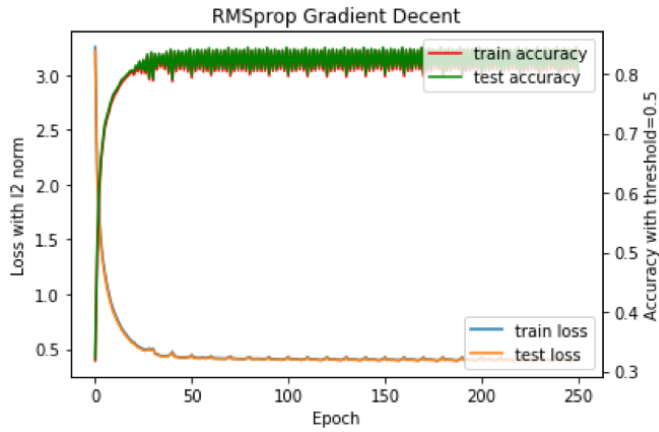


Fig. 15. Metrics evolution process details

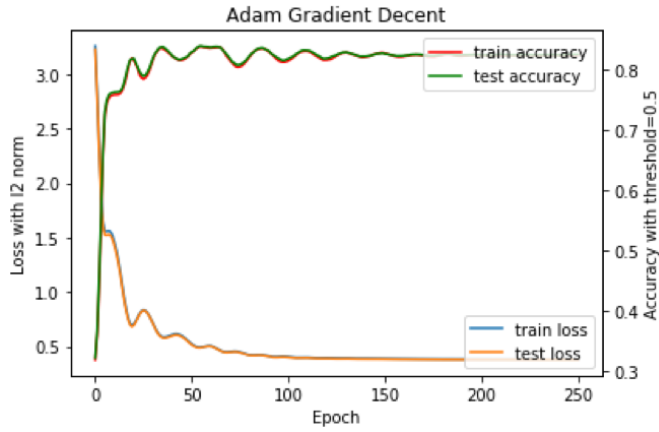


Fig. 16. Metrics evolution process details

algorithms. *Neural Networks; The Official Journal of the International Neural Network Society*, 12(1), 145151.

[4] Nesterov, Y. (1983). A method for unconstrained convex minimization problem with the rate of convergence  $o(1/k^2)$ . *Doklady ANSSSR* (translated as *Soviet.Math.Docl.*), vol. 269, pp. 543 547.

[5] Bengio, Y., Boulanger-Lewandowski, N., & Pascanu, R. (2012). *Advances in Optimizing Recurrent Networks*. Retrieved from <http://arxiv.org/abs/1212.0901>

[6] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method. Retrieved from <http://arxiv.org/abs/1212.5701>

[7] Kingma, D. P., & Ba, J. L. (2015). Adam: a Method for Stochastic Optimization. *International Conference on Learning Representations*, 113.

REFERENCES

[1] Ruder S. An overview of gradient descent optimization algorithms[J]. 2016.

[2] Sutton, R. S. (1986). Two problems with backpropagation and other steepest-descent learning procedures for networks. *Proc. 8th Annual Conf. Cognitive Science Society*.

[3] Qian, N. (1999). On the momentum term in gradient descent learning