



华南理工大学

South China University of Technology

The Experiment Report of *Machine Learning*

College	<u>Software College</u>
Subject	<u>Software Engineering</u>
Members	<u>JunPeng Su</u>
Student ID	<u>201530612743</u>
E-mail	<u>sujunpeng.chn@gmail.com</u>
Tutor	<u>MingKui Tan</u>
Date submitted	<u>2017. 12 . 14</u>

1. Topic:

Comparison of Various Stochastic Gradient Descent Methods for Solving Classification Problems

2. Time: 2017-12-02 2:00-5:00 PM B7-138

3. Reporter: JunPeng Su

4. Purposes:

- 1) Compare and understand the difference between Gradient Descent and Stochastic Gradient Descent.
- 2) Compare and understand the difference and relationship between Logistic Regression and Linear Classification.
- 3) Further understand the principles of SVM and practice on larger data.
- 4) Compare the performance of different variants of Stochastic Gradient Descent

5. Data sets and data analysis:

Experimental dataset is a9a of LIBSVM Data, including 32561/16281(testing) samples and each sample has 123/123 (testing) features.

6. Experimental steps:

Logistic Regression

- 1) Load train data and validation data.
- 2) Set the hyper-parameters of Logistic Regression model and select an appropriate threshold.
- 3) Initialize linear regression model parameters by setting all parameters into 0.
- 4) Choose the loss function and calculate the derivative.
- 5) Calculate gradient towards the loss function from partial samples
- 6) Update linear regression model using different optimization methods (NAG, RMSProp, AdaDelta, Adam).
- 7) Mark samples whose predict score greater than threshold as positive , on the contrary as negative.
- 8) Predict under validation dataset and calculate the loss
- 9) Repeat step 5) - 8) several times
- 10) Draw graph of loss with the number of iterations

Linear Classification

- 1) Load train data and validation data.
- 2) Set the hyper-parameters of Linear Classification model and select an appropriate threshold.

- 3) Initialize SVM model parameters by setting all parameters into 0.
- 4) Choose the loss function and calculate the derivative.
- 5) Calculate gradient towards the loss function from partial samples
- 6) Update linear regression model using different optimization methods (NAG, RMSProp, AdaDelta, Adam).
- 7) Mark samples whose predict score greater than threshold as positive , on the contrary as negative.
- 8) Predict under validation dataset and calculate the loss
- 9) Repeat step 5) - 8) several times
- 10) Draw graph of loss with the number of iterations

7. Code:

Logistic Regression

#calculate gradient and loss

```

14 #define loss function
15 def loss(w,x,y):
16     #get data volume
17     data_volume=x.shape[0]
18     #calculate positive probability
19     positive_probability=sigmoid(w,x)
20     #calculate negative probability using positive probability
21     negative_probability=1-positive_probability
22     #sum the loss
23     loss_total=np.dot(y.T,np.log(positive_probability))+np.dot((1-y).T,np.log(negative_probability))
24     #mean the loss
25     result=-loss_total/data_volume
26     return result
27
28 #define gradient calculating function
29 def gradient(w,x,y):
30     #get data volume
31     data_volume=x.shape[0]
32     #calculate positive probability
33     positive_probability=sigmoid(w,x)
34     #calculate gradient
35     grad=np.dot((positive_probability-y).T,x).T/data_volume
36     return grad

```

#RMSProp optimization method

```

81 #define RMSProp optimization method
82 def RMSProp_Optimization(threshold,w,x_train,y_train,x_test,y_test,runs):
83     #print a signal
84     print('RMSProp Optimization')
85     #set parameter
86     learning_rate=0.02
87     gama=0.9
88     e=0.0000001
89     #initialize medium variable
90     G=0
91     #initialize container for loss
92     test_losses=np.zeros(runs)
93     #optimize
94     for i in range(runs):
95         #generate random integer
96         index=np.arange(5000)
97         np.random.shuffle(index)
98         #get a batch of stochastic training sample
99         x_train_part=x_train[index]
100        y_train_part=y_train[index]
101        #calculate gradient
102        grad=gradient(w,x_train_part,y_train_part)
103        G=gama*G+(1-gama)*np.dot(grad.T,grad)
104        w=w-learning_rate/np.sqrt(G+e)*grad
105        #calculate the loss
106        error=loss(w,x_test,y_test)
107        #save the loss in container
108        test_losses[i]=error
109        #Log every 100 epochs
110        if(i%100==99):
111            #calculate the accuracy
112            accur=accuracy(threshold,w,x_test,y_test)
113            print("epochs "+str(i+1)+" Accuracy: "+str(accur))
114        #draw graph of loss
115        x_axis=np.arange(runs)
116        plt.plot(x_axis,test_losses,'r',label='RMSProp')
117        return w

```

#preprocess data

```
290 #Load data
291 train_data=load_svmlight_file("a9a.txt")
292 test_data=load_svmlight_file("a9a.t",n_features=123)
293 #separate features and labels
294 X_train,y_train=train_data[0],train_data[1]
295 X_test,y_test=test_data[0],test_data[1]
296 #change the range of y from [-1,1] to [0,1]
297 y_train=y_train/2+0.5
298 y_test=y_test/2+0.5
299 #change sparse matrix to numpy array
300 X_train=X_train.toarray()
301 X_test=X_test.toarray()
302 #integrate features with bias
303 #create bias vector
304 X_train_bias=np.ones((X_train.shape[0],1))
305 X_test_bias=np.ones((X_test.shape[0],1))
306 #add bias column to features
307 X_train=np.hstack([X_train,X_train_bias])
308 X_test=np.hstack([X_test,X_test_bias])
```

Linear Classification

#calculate the gradient and loss

```
13 #define loss function
14 def loss(C,threshold,w,x,y):
15     loss_part_1=np.dot(w.T,w)/2
16     loss_part_2=0
17     #get the volume of data
18     data_volume=x.shape[0]
19     #calculate loss for each record
20     for i in range(data_volume):
21         hinge=1-y[i]*(np.dot(x[i],w))
22         if hinge>=0:
23             loss_part_2+=hinge
24     #calculate the total loss
25     loss_all=(loss_part_1+C*loss_part_2)/data_volume
26     return loss_all
27
28 #define gradient calculating function
29 def gradient(C,threshold,w,b,x,y):
30     grad_w=np.zeros_like(w)
31     grad_b=0
32     #get the volume of data
33     data_volume=x.shape[0]
34     #calculate gradient for each record
35     for i in range(data_volume):
36         hinge=1-y[i]*(np.dot(x[i],w))
37         #deal with uncontinuous gradient
38         if hinge>=threshold:
39             grad_w=grad_w-C*y[i]*x[i].T
40             grad_b=grad_b-C*y[i]
41
42     grad_w=grad_w/w
43     return grad_w/data_volume,grad_b/data_volume
44
```

#NAG optimization method

```
230 #define NAG optimization method
231 def NAG_Optimization(C,threshold,w,b,x_train,y_train,x_test,y_test,runs):
232     #print a signal
233     print('NAG')
234     #set parameter
235     gama=0.9
236     learning_rate=0.01
237     #initialize medium variables
238     v_w=np.zeros_like(w)
239     v_b=0
240     #initialize container for loss
241     test_losses=np.zeros(runs)
242     #optimize
243     for i in range(runs):
244         #generate random integer
245         index=np.arange(5000)
246         np.random.shuffle(index)
247         #get a batch of stochastic training sample
248         x_train_part=x_train[index]
249         y_train_part=y_train[index]
250         #calculate gradient
251         grad_w,grad_b=gradient(C,threshold,w,b,x_train_part,y_train_part)
252         grad_w=grad_w+gama*v_w
253         grad_b=grad_b+gama*v_b
254         #update v
255         v_w=gama*v_w+learning_rate*grad_w
256         v_b=gama*v_b+learning_rate*grad_b
257         #update weight vector
258         w=w-v_w
259         b=b-v_b
260         #calculate the loss
261         error=loss(C,threshold,w,x_test,y_test)
262         #save the loss in container
263         test_losses[i]=error
264         #Log every 100 epochs
265         if(i%100==99):
266             #calculate the accuracy
267             accur=accuracy(threshold,w,b,x_test,y_test)
268             print("epochs "+str(i+1)+" Accuracy: "+str(accur))
269     #draw graph of Loss
270     x_axis=np.arange(runs)
271     plt.plot(x_axis,test_losses,'b',label='NAG')
272     return w
```

#preprocess the data

```
314 #Load data
315 train_data=load_svmlight_file("a9a.txt")
316 test_data=load_svmlight_file("a9a.t",n_features=123)
317 #separate features and labels
318 X_train,y_train=train_data[0],train_data[1]
319 X_test,y_test=test_data[0],test_data[1]
320 #change sparse matrix to numpy array
321 X_train=X_train.toarray()
322 X_test=X_test.toarray()
```

p.s. Only part of the codes are displayed here, the whole codes are in RegressionExperiment.ipynb and ClassificationExperiment.ipynb.

8. The initialization method of model parameters:

Set all model parameters into 0.

9. The selected loss function and its derivatives:

Logistic Regression

Loss function:

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

Derivatives:

$$\frac{\partial J(w)}{\partial w} = (h_w(X) - y)X$$

Linear Classification

loss function:

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

derivatives:

$$\frac{\partial L_D(w, b)}{\partial w} = w + C \sum_{i=1}^n g_w(x_i)$$

$$\frac{\partial L_D(w, b)}{\partial b} = C \sum_{i=1}^n g_b(x_i)$$

here,

$$g_w(x_i) = \begin{cases} -y_i x_i & 1 - y_i(w^T x_i + b) \geq 0 \\ 0 & 1 - y_i(w^T x_i + b) < 0 \end{cases}$$

$$g_b(x_i) = \begin{cases} -y_i & 1 - y_i(w^T x_i + b) \geq 0 \\ 0 & 1 - y_i(w^T x_i + b) < 0 \end{cases}$$

10. Experimental results and curve:

Hyper-parameter selection:

Logistic Regression

RMSProp			AdaDelta		Adam				NAG		Momentum	
η	γ	e	γ	e	β	γ	η	e	γ	η	γ	η
0.001	0.9	1e-7	0.95	1e-7	0.9	0.999	0.1	1e-7	0.9	0.001	0.9	0.0001
0.005	0.9	1e-7	0.90	1e-7	0.9	0.999	0.05	1e-7	0.9	0.005	0.9	0.001
0.01	0.9	1e-7	0.60	1e-7	0.8	0.999	0.1	1e-7	0.9	0.01	0.9	0.05
0.01	0.9	1e-7	0.60	1e-5	0.8	0.999	0.1	1e-7	0.9	0.01	0.9	0.05
0.1	0.9	1e-7	0.80	1e-5	0.8	0.999	0.1	1e-6	0.9	0.1	0.9	0.05
0.05	0.9	1e-7	0.80	1e-4	0.8	0.8	0.1	1e-6	0.9	0.05	0.9	0.05
0.02	0.9	1e-7	0.80	1e-4	0.8	0.9	0.1	1e-6	0.9	0.01	0.9	0.05

p.s. To compare the performance of different optimization methods, epochs used in all methods are set into the same number 1000.

Linear Classification

RMSProp			AdaDelta		Adam				NAG		Momentum	
η	γ	e	γ	e	β	γ	η	e	γ	η	γ	η
0.005	0.9	1e-7	0.95	1e-7	0.9	0.999	0.1	1e-7	0.9	0.001	0.9	0.01
0.02	0.9	1e-7	0.80	1e-4	0.8	0.9	0.1	1e-6	0.9	0.01	0.9	0.05
0.1	0.9	1e-7	0.95	1e-5	0.9	0.999	0.2	1e-7	0.9	0.001	0.9	0.001
0.005	0.9	1e-7	0.95	1e-5	0.9	0.999	0.05	1e-7	0.9	0.1	0.9	0.1
0.005	0.9	1e-7	0.95	1e-5	0.9	0.999	0.05	1e-7	0.9	0.05	0.9	0.08
0.005	0.9	1e-7	0.95	1e-5	0.9	0.999	0.05	1e-7	0.9	0.01	0.9	0.02
0.005	0.9	1e-7	0.95	1e-5	0.3	0.999	0.2	1e-7	0.9	0.01	0.9	0.02

p.s. To compare the performance of different optimization methods, epochs used in all methods are set into the same number 400.

Predicted Results (Best Results):

Logistic Regression

RMSProp			AdaDelta		Adam				NAG		Momentum	
η	γ	e	γ	e	β	γ	η	e	γ	η	γ	η
0.02	0.9	1e-7	0.80	1e-4	0.8	0.9	0.1	1e-6	0.9	0.01	0.9	0.05

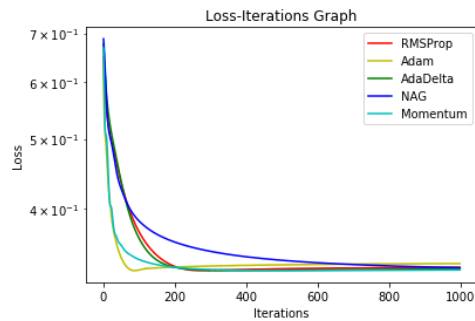
Linear Classification

RMSProp			AdaDelta		Adam				NAG		Momentum	
η	γ	e	γ	e	β	γ	η	e	γ	η	γ	η

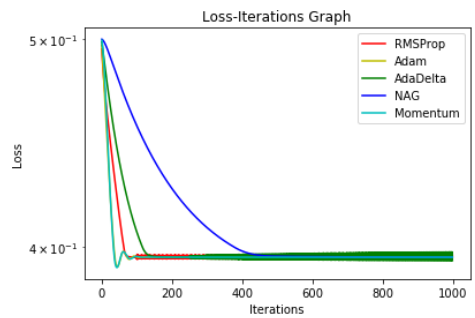
0.005	0.9	1e-7	0.95	1e-5	0.3	0.999	0.2	1e-7	0.9	0.01	0.9	0.02
-------	-----	------	------	------	-----	-------	-----	------	-----	------	-----	------

Loss curve:

Logistic Regression



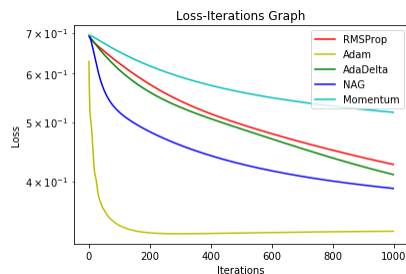
Linear Classification



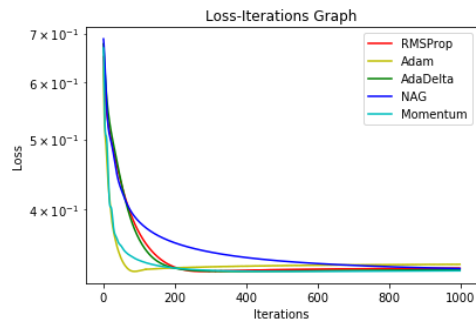
11. Results analysis:

Logistic Regression

The loss curves below are generated in the first try when all hyper-parameters are choose randomly. It indicates that Adam Optimization method performs best. As an improvement version of Momentum Optimization method, NAG Optimization method performs much better than Momentum. And RMSProp Optimization method and AdaDelta Optimization method are comparable.



The next loss curves are generated after the process of tuning hyper-parameters. It shows that selecting appropriate hyper-parameters improve the performance and all optimization methods are able to obtain similar test loss at last.

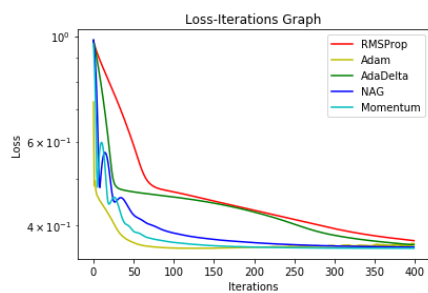


Accuracy of different optimization method is shown in the table below. NAG obtain the highest accuracy while Adam obtain the lowest.

Optimization Method	RMSProp	AdaDelta	Adam	NAG	Momentum
Accuracy	0.84381	0.84467	0.84264	0.84860	0.84546

Linear Classification

The loss curves below are generated after the process of tuning hyper-parameters. The graph shows that Adam optimization method and Momentum optimization method have the best performance. And all method obtain similar test loss at last.



Accuracy of different optimization method is shown in the table below. Momentum obtain the highest accuracy while RMSProp obtain the lowest.

Optimization Method	RMSProp	AdaDelta	Adam	NAG	Momentum
Accuracy	0.80075	0.81316	0.81795	0.84215	0.84706

Comparison between loss curves of logistic regression and linear classification indicates that Adam optimization method is the fastest to reduce the test loss. However, accuracy of the model optimized by Adam optimization method is unsatisfied which means that using Adam optimization method is prone to over-fitting.

12. Similarities and differences between logistic regression and linear classification :

Similarities: Logistic regression and linear classification are both attempting to represent the relationship between features and labels in a linear way. They both use the formula $y=Wx$.

Difference: Logistic regression need to transform y to a real value between 0 and 1 which represents the probability of positive or negative sample.

13. Summary:

The loss curves show that Adam optimization method is the fastest optimization method among all methods. Adam optimization method reduces the loss fast and is easy to select hyper-parameters while it leads to over-fitting sometimes. From the accuracy table, I learned that NAG optimization method and Momentum optimization method perform well.

By comparing the loss curves of the same optimization method using different hyper-parameters, I realize the importance of selecting appropriate hyper-parameters. After a hard struggle in tuning hyper-parameters, I learned some skills about selecting appropriate hyper-parameters. Too large learning rate will make the loss curve fluctuates heavily. Different hyper-parameters have different influence on the performance such as hyper-parameter ϵ behaves better than γ in speeding up the process of loss reduction.