



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:  
Guobin Wu

Supervisor:  
Qingyao Wu

Student ID:  
201530613030

Grade:  
Undergraduate

December 11, 2017

# Logistic Regression, Linear Classification and Stochastic Gradient Descent

**Abstract—This report tends to illustrate the experiments we have done about logistic regression, linear classification and stochastic gradient descent(SGD), with respect to understanding and comprehending the core of this mentioned topics.**

## I. INTRODUCTION

Logistic regression and linear classification are the two of most fundamental machine learning models. Additionally, gradient decent(GD) is one of the most widely-used optimizing methods to reach local optimal solution. Stochastic gradient descent(SGD), an improved version of traditional GD, accelerates the process reaching the solution. This experiment aims to compare GD to SGD, to help understanding the differences and relations between them. What's more, we also compare logistic regression to linear classification, figuring out what is and is not similar to each other. Lastly, we practice SVM on larger data to have a better command of its principles.

## II. METHODS AND THEORY

### *Logistic Regression and Stochastic Gradient Descent*

1. Load the training set and validation set.
2. Initialize logistic regression model parameters, you can consider initializing zeros, random numbers or normal distribution.
3. Select the loss function and calculate its derivation, find more detail in PPT.
4. Calculate gradient G 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.

---

### *Linear Classification and Stochastic Gradient Descent*

1. Load the training set and validation set.
  2. Initialize SVM model parameters, you can consider initializing zeros, random numbers or normal distribution.
  3. Select the loss function and calculate its derivation, find more detail in PPT.
  4. Calculate gradient G 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.

## III. EXPERIMENT

### A. Data Set

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

### B. Implementation

#### 1) Logistic regression:

Loss function of logistic regression is as follows.

$$L(\theta) = -\frac{1}{m} \sum_{i=0}^n (y_i \log(h_\theta(X_i)) + (1 - y_i) \log(1 - h_\theta(X_i)))$$

$$h_\theta(X) = \frac{1}{1+e^{-\theta^T \cdot X}}$$

where

Compute the gradient of  $L(\theta)$ , we obtain,

$$\frac{\partial L}{\partial \theta} = \frac{1}{n} \sum_{i=0}^n X_i (h_\theta(X) - y_i)$$

Having defined loss function and its gradient, we can use SGD to get the final solution. We use four method respectively to reach the local optimal solution including NAG, RMSProp, AdaDelta and Adam. The super parameters we select are as follows.

NAG	learning rate	0.005
	gamma	0.9
	iteration	1000
RMSProp	learning rate	0.005
	gamma	0.9
	epsilon	1e-8
	iteration	1000
AdaDelta	learning rate	0.005
	gamma	0.9

	epsilon	1e-8
	iteration	1000
Adam	learning rate	0.005
	beta1	0.9
	beta2	0.999
	epsilon	1e-8
	iteration	1000

Next, we program to implement the above methods. The following are the screenshots of source code.

```

1  from sklearn import datasets, model_selection, linear_model
2  import numpy as np
3  import jupyter
4  import matplotlib.pyplot as plt
5  import math
6  import random
7
8  X_train, y_train = datasets.load_svmlight_file("a9atrain.txt")
9
10 # turn the csr_matrix into array for futher processing
11 x_ = np.array(X_train.toarray(), np.float32).reshape((-1, 123))
12 y_ = np.array(y_train, np.float32).reshape((-1, 1))
13
14 for i in range(y_.shape[0]):
15     if y_[i,0] == -1.0 : y_[i,0] = 0.
16
17 X_ = np.hstack([x_, np.ones((x_.shape[0], 1))])
18
19
20 X_test, y_test = datasets.load_svmlight_file("a9atest.txt")
21
22 xt_ = np.array(X_test.toarray(), np.float32).reshape((-1, 122))
23 yt_ = np.array(y_test, np.float32).reshape((-1, 1))
24 for i in range(yt_.shape[0]):
25     if yt_[i,0] == -1.0 : yt_[i,0] = 0.
26 # X_ = (x_-1) to fit the constant item
27
28 Xt_ = np.hstack([xt_, np.zeros((xt_.shape[0], 1)),np.ones((xt_.shape[0], 1))])
29
30
31 # hθ(X) = e^(Theta * X) / (1 + e^(Theta * X)) = 1 / (1 + e^(- Theta *
32 def h_theta(Xi, Theta):
33     e_t = math.exp(Xi.dot(Theta.T))
34     return e_t / (1 + e_t)
35
36 # L(θ) = - (1/m) ∑ (yi * log(hθ(Xi)) + (1 - yi) * log(1 - hθ(Xi)))
37 def compute_loss(X, y, Theta):
38     m = y.shape[0]
39     loss = 0.
40     for i in range(m):
41         loss += (y[i] * math.log(h_theta(X[i, :], Theta))) + ((1 - y[i]) * math.log(1 - h_theta(X[i, :], Theta)))
42     loss /= - m
43     return loss

```

```

44
45 # ∂L/∂θ = (1/m) ∑ (hθ(Xi) - yi) * Xi
46 def compute_gradient(X, y, Theta):
47     m = y.shape[0]
48     gradient = np.zeros(Theta.shape)
49     for i in range(m):
50         gradient += (h_theta(X[i, :], Theta) - y[i]) * (X[i, :])
51     gradient /= m
52     return gradient
53
54
55 def train_model_nag(X, y, Theta, learning_rate, gamma, iteration = 10):
56     test_loss_history = np.zeros((iteration, 1))
57     v = np.zeros(Theta.shape)
58     Theta_gradient = np.zeros(Theta.shape)
59     for iter in range(iteration):
60         index = random.randint(0, y.shape[0]-10)
61         Theta = Theta - gamma * v
62         v = gamma * v - learning_rate * compute_gradient(X[index:index+10,:], y[index:index+10])
63         Theta = Theta + v
64         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
65     return test_loss_history, Theta
66
67
68 def train_model_rmsprop(X, y, Theta, learning_rate, gamma, epsilon, iteration = 10):
69     test_loss_history = np.zeros((iteration, 1))
70     G_t = 0.
71     Theta_gradient = np.zeros(Theta.shape)
72     for iter in range(iteration):
73         index = random.randint(0, y.shape[0]-10)
74         Theta_gradient = compute_gradient(X[index:index+10,:], y[index:index+10])
75         G_t = gamma * G_t + (1 - gamma) * Theta_gradient.dot(Theta_gradient)
76         Theta = Theta - (learning_rate / np.sqrt(G_t + epsilon)) * Theta_gradient
77         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
78     return test_loss_history, Theta
79
80
81 def train_model_adadelta(X, y, Theta, gamma, epsilon, iteration = 10):
82     test_loss_history = np.zeros((iteration, 1))
83     Theta_gradient = np.zeros(Theta.shape)
84     G_t = 0.
85     delta_theta = np.zeros(Theta.shape)
86     delta_t = 0.03
87     for iter in range(iteration):
88         index = random.randint(0, y.shape[0]-10)
89         Theta_gradient = compute_gradient(X[index:index+10,:], y[index:index+10])
90         G_t = gamma * G_t + (1 - gamma) * Theta_gradient.dot(Theta_gradient)
91         delta_theta = - (np.sqrt(delta_t + epsilon) / np.sqrt(G_t + epsilon)) * Theta_gradient
92         Theta = Theta + delta_theta
93         delta_t = gamma * delta_t + (1 - gamma) * (delta_theta * delta_theta)
94         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)
95     return test_loss_history, Theta
96
97 def train_model_adam(X, y, Theta, learning_rate, beta1, beta2, iteration = 1000):
98     test_loss_history = np.zeros((iteration, 1))
99     Theta_gradient = np.zeros(Theta.shape)
100    v_t = 0.
101    m_t = np.zeros(Theta.shape)
102    for iter in range(iteration):
103        index = random.randint(0, y.shape[0]-10)
104        Theta_gradient = compute_gradient(X[index:index+10,:], y[index:index+10])
105        m_t = beta1 * m_t + (1 - beta1) * Theta_gradient
106        v_t = beta2 * v_t + (1 - beta2) * Theta_gradient.dot(Theta_gradient)
107        mt_estimate = m_t / (1 - pow(beta1, iter + 1))
108        vt_estimate = v_t / (1 - pow(beta2, iter + 1))
109        Theta = Theta - learning_rate * mt_estimate / (np.sqrt(vt_estimate + epsilon))
110        test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)
111    return test_loss_history, Theta
112
113 iteration = 1000
114
115 beta1 = 0.9
116 beta2 = 0.999
117 epsilon = 1e-8

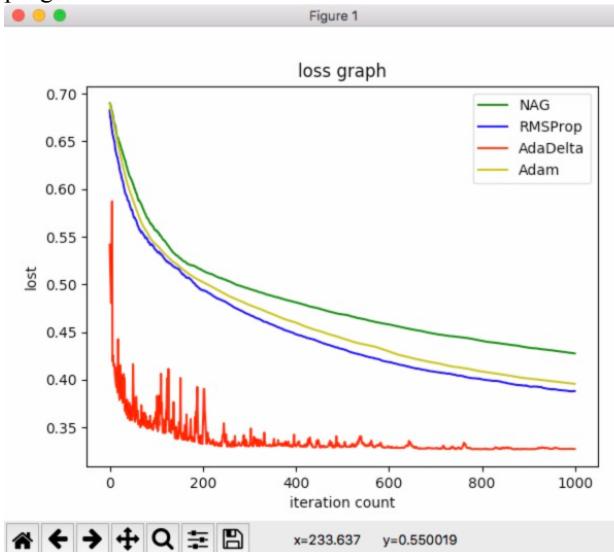
```

```

121 t_nag = np.zeros((1, 124))
122 t_rmsprop = np.zeros((1, 124))
123 t_adadelta = np.zeros((1, 124))
124 t_adam = np.zeros((1, 124))
125 # for i in range(t.shape[0]):
126 | # t[i] = [-0.03]
127
128 nag_loss_history, t_nag = train_model_nag(X, y, t_nag, 0.001)
129 rmsprop_loss_history, t_rmsprop = train_model_rmsprop(X, y, t_rmsprop, 0.001)
130 adadelta_loss_history, t_adadelta = train_model_adadelta(X, y, t_adadelta, 0.001)
131 adam_loss_history, t_adam = train_model_adam(X, y, t_adam, 0.001)
132
133
134 plt.plot(nag_loss_history, 'g', label='NAG')
135 plt.plot(rmsprop_loss_history, 'b', label='RMSProp')
136 plt.plot(adadelta_loss_history, 'r', label='AdaDelta')
137 plt.plot(adam_loss_history, 'y', label='Adam')
138
139
140 plt.legend(loc='upper right')
141
142 plt.ylabel('loss');
143
144 plt.xlabel('iteration count')
145
146 plt.title('loss graph')
147
148 plt.show()

```

We get the following loss graphs as results after running the program.



From the graph we find that AdaDelta reaches the local optimal solution fastest, but it also has obvious vibration. By contrast, NAG is slower than AdaDelta with a smoother curve. RMSProp and Adam are closely overlapped, which are slower than AdaDelta and faster than NAG. Four methods reaches optimal solution far faster than traditional GD.

## 2) Linear classification:

Loss function of logistic regression is as follows.

$$L(\theta) = \frac{1}{2n} \sum_{i=0}^n (y_i - h_\theta(X_i))^2$$

, where  $h_\theta(X) = \sum_{i=0}^n \theta_i X_i$

Compute the gradient of  $L(\theta)$ , we obtain,

$$\frac{\partial L}{\partial \theta} = \frac{1}{n} \sum_{i=0}^n (y_i - h_\theta(X_i)) \cdot X_i$$

Having defined loss function and its gradient, we can use SGD to get the final solution. We use four method respectively to reach the local optimal solution including NAG, RMSProp, AdaDelta and Adam. The super parameters we select are as follows.

NAG	learning rate	0.005
	gamma	0.9
	iteration	3000
RMSProp	learning rate	0.005
	gamma	0.9
	epsilon	1e-8
	iteration	3000
AdaDelta	learning rate	0.005
	gamma	0.9
	epsilon	1e-8
	iteration	3000
Adam	learning rate	0.005
	beta1	0.9
	beta2	0.999
	epsilon	1e-8
	iteration	3000

Next, we program to implement the above methods. The following are the screenshots of source code.

```

1 import numpy
2 import random
3 import jupyter
4 import math
5 from sklearn.datasets import load_svmlight_file
6 from sklearn.model_selection import train_test_split
7 from matplotlib import pyplot
8
9 x, y_train = load_svmlight_file("a9atrain.txt")
10 x_train = x.toarray()
11 x, y_test = load_svmlight_file("a9atest.txt")
12 x_test = x.toarray()
13
14 X_train = numpy.hstack([x_train, numpy.ones((x_train.shape[0], 1))])
15 X_test = numpy.hstack([x_test, numpy.zeros((x_test.shape[0], 1))])
16 X_test = numpy.hstack([X_test, numpy.ones((x_test.shape[0], 1))])
17
18 def compute_grad(x, y, w):
19     gradient = x * (y - x.dot(w.T))
20     return gradient
21
22 def compute_loss(x, y, w, random_i):
23     loss = 0
24     a = len(random_i)
25     for m in range(a):
26         loss += 0.5 * ((y[random_i[m]] - x[random_i[m]].dot(w)))
27     return loss/a
28

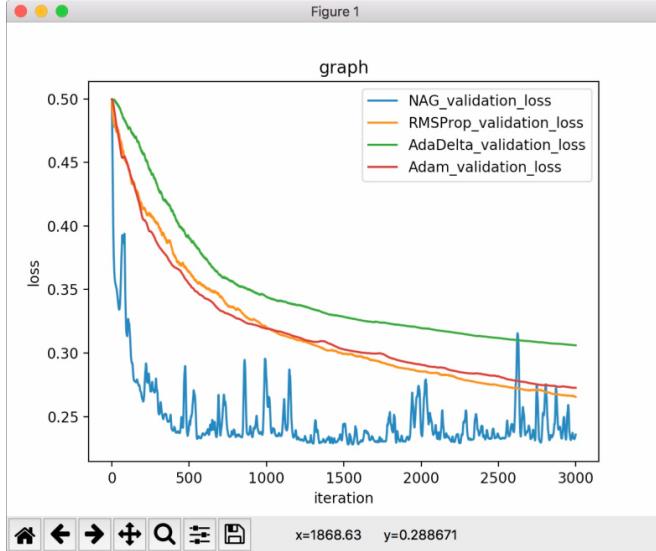
```

```

29 def NAG_train(x, y, x_test, y_test, w, C, lr, gamma,
30     vt = numpy.zeros(w.shape)
31     loss_history = []
32     test_loss_history = []
33     random_index = []
34     random_test_index = []
35     for i in range(iteration):
36         random_num = random.randint(0, x.shape[0]-1)
37         random_test_num = random.randint(0, x_test.shape[0]-1)
38         random_index.append(random_num)
39         random_test_index.append(random_test_num)
40     for i in range(iteration):
41         gradient = compute_grad(x[random_index[i],:])
42         vt = gamma*vt - lr*gradient
43         w -= vt
44         loss = compute_loss(x, y, w, random_index)
45         loss_history.append(loss)
46         test_loss_history.append(compute_loss(x_test, y_test))
47         if loss < threshold :
48             break
49     return w, loss_history, test_loss_history
50
51 def RMSProp_train(x, y, x_test, y_test, w, C, lr, gamma,
52     Gt = 0
53     loss_history = []
54     test_loss_history = []
55     random_index = []
56     random_test_index = []
57     for i in range(iteration):
58         random_num = random.randint(0, x.shape[0]-1)
59         random_test_num = random.randint(0, x_test.shape[0]-1)
60         random_index.append(random_num)
61         random_test_index.append(random_test_num)
62     for i in range(iteration):
63         gradient = compute_grad(x[random_index[i],:])
64         Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient)
65         w += lr * gradient / math.sqrt(Gt+1e-8)
66         loss = compute_loss(x, y, w, random_index)
67         loss_history.append(loss)
68         test_loss_history.append(compute_loss(x_test, y_test))
69         if loss < threshold :
70             break
71     return w, loss_history, test_loss_history
72
73 def AdaDelta_train(x, y, x_test, y_test, w, C, lr, gamma,
74     Gt = 0
75     variable_t = 0
76     loss_history = []
77     test_loss_history = []
78     random_index = []
79     random_test_index = []
80     for i in range(iteration):
81         random_num = random.randint(0, x.shape[0]-1)
82         random_test_num = random.randint(0, x_test.shape[0]-1)
83         random_index.append(random_num)
84         random_test_index.append(random_test_num)
85     for i in range(iteration):
86         gradient = compute_grad(x[random_index[i],:], y[random_index[i]])
87         Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient)
88         variable_w = - math.sqrt(variable_t + 1e-8) * gradient
89         w -= variable_w
90         variable_t = gamma*variable_t + (1-gamma)*variable_w**2
91         loss = compute_loss(x, y, w, random_index)
92         loss_history.append(loss)
93         test_loss_history.append(compute_loss(x_test, y_test))
94         if loss < threshold :
95             break
96     return w, loss_history, test_loss_history
97
98 def Adam_train(x, y, x_test, y_test, w, C, lr, gamma, threshold):
99     Gt = 0
100    moment = numpy.zeros((1, x.shape[1]))
101    B = 0.9
102    loss_history = []
103    test_loss_history = []
104    random_index = []
105    random_test_index = []
106    for i in range(iteration):
107        random_num = random.randint(0, x.shape[0]-1)
108        random_test_num = random.randint(0, x_test.shape[0]-1)
109        random_index.append(random_num)
110        random_test_index.append(random_test_num)
111    for i in range(iteration):
112        gradient = compute_grad(x[random_index[i],:], y[random_index[i]])
113        moment = B*moment + (1-B)*gradient
114        Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient.T)
115        a = lr * math.sqrt(1 - pow(gamma, iteration)) / (1 - pow(B, iteration))
116        w += a * moment / math.sqrt(Gt + 1e-8)
117        loss = compute_loss(x, y, w, random_index)
118        loss_history.append(loss)
119        test_loss_history.append(compute_loss(x_test, y_test))
120        if loss < threshold :
121            break
122    return w, loss_history, test_loss_history
123
124 iteration = 3000
125 # NAG
126 NAG_w = numpy.zeros((1, X_train.shape[1]))
127 NAG_w, NAG_loss_history, NAG_test_loss_history = NAG_train(X_train, y_train, X_test, y_test, NAG_w, iteration)
128
129 # RMSProp
130 RMS_w = numpy.zeros((1, X_train.shape[1]))
131 RMS_w, RMS_loss_history, RMS_test_loss_history = RMSProp_train(X_train, y_train, X_test, y_test, RMS_w, iteration)
132
133 # AdaDelta
134 AdaDelta_w = numpy.zeros((1, X_train.shape[1]))
135 AdaDelta_w, AdaDelta_loss_history, AdaDelta_test_loss_history = AdaDelta_train(X_train, y_train, X_test, y_test, AdaDelta_w, iteration)
136
137 #Adam
138 Adam_w = numpy.zeros((1, X_train.shape[1]))
139 Adam_w, Adam_loss_history, Adam_test_loss_history = Adam_train(X_train, y_train, X_test, y_test, Adam_w, iteration)
140
141
142 pyplot.plot(NAG_test_loss_history, label = 'NAG_validation')
143 pyplot.plot(RMS_test_loss_history, label = 'RMSProp_validation')
144 pyplot.plot(AdaDelta_test_loss_history, label = 'AdaDelta_validation')
145 pyplot.plot(Adam_test_loss_history, label = 'Adam_validation')
146 pyplot.legend(loc='upper right')
147 pyplot.ylabel('loss')
148 pyplot.xlabel('iteration')
149 pyplot.title('graph')
150 pyplot.show()

```

We get the following loss graphs as results after running the program.



From the graph we find that NAG reaches the local optimal solution fastest, but it also has obvious vibration. By contrast, AdaDelta is slower than NAG with a smoother curve. RMSProp and Adam are closely overlapped, which are slower than NAG and faster than AdaDelta. Four methods reaches optimal solution far faster than traditional GD.

#### IV. CONCLUSION

SGD improves traditional GD with a faster converging speed and considerable result. Four specific methods are different improvement of plain SGD with respective advantages. Logistic regression and linear classification both solves classification problem to predict new samples. The major difference between them is the way they evaluate final super plane. Logistic regression tries to find a super plane which makes all sample points get away from it, whereas linear classification finds a plane defined by so-called ‘support vectors’. This experiment makes comparisons between GD and SGD, logistic regression and linear classification, evaluates four SGD methods in this two model and analyzes the experiment results.