# The Experiment Report of
# *Machine Learning*

**College**  Software College

**Subject**  Software Engineering

**Members**

**Student ID**  201530612651

**E-mail**  835380624@qq.com

**Tutor**  mingkui Tan

**Date submitted**  2017.12. 11

1. **Topic:**

   Comparison of Various Stochastic Gradient Descent Methods for

   Solving Classification Problems

2. **Time:**

   2017-12-02 2:00-5:00 PM

3. **Reporter:**

   任嘉宁

4. **Purposes:**

   ● 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.

5. **Data sets and data analysis:**

   Experiment uses a9a of LIBSVM Data, including

   32561/16281(testing) samples and each sample has 123/123

   (testing) features.

6. **Experimental steps:**

   Logistic Regression and Stochastic Gradient Descent：

   ● Load the training set and validation set.

   ● Initialize logistic regression model parameters with random

numbers.

- Select the loss function and calculate its derivation.

- Calculate gradient $\mathcal{G}$ 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 $\mathcal{L}_{NAG}$, $\mathcal{L}_{RMSProp}$, $\mathcal{L}_{AdaDelta}$ and $\mathcal{L}_{Adam}$.

- Repeat step 4 to 6 for several times, and drawing graph of $\mathcal{L}_{NAG}$, $\mathcal{L}_{RMSProp}$, $\mathcal{L}_{AdaDelta}$ and $\mathcal{L}_{Adam}$ with the number of iterations.

Linear Classification and Stochastic Gradient Descent:

- Load the training set and validation set.

- Initialize SVM model parameters with random numbers.

- Select the loss function and calculate its derivation.

- Calculate gradient $\mathcal{G}$ 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 $\mathcal{L}_{NAG}$，$\mathcal{L}_{RMSProp}$，$\mathcal{L}_{AdaDelta}$ and $\mathcal{L}_{Adam}$.

- Repeat step 4 to 6 for several times, and drawing graph of $\mathcal{L}_{NAG}$，$\mathcal{L}_{RMSProp}$，$\mathcal{L}_{AdaDelta}$ and $\mathcal{L}_{Adam}$ with the number of iterations.

## 7. Code:

Logistic Regression and Stochastic Gradient Descent：

NAG:

```python
# NAG

learning_rate=0.05
w=0.001*np.random.random((1,124))
v=np.zeros((1,124))
gamma=0.9

loss_NAGhistory=[]
test_NAGhistory=[]
for k in range(time):
    mask=np.random.choice(X_train.shape[0],256,replace=False)
    X_batch=X_train[mask]
    y_batch=y_train[mask]
    NAGloss=np.mean(np.log(1+np.exp(-y_batch.reshape(-1,1)*(X_batch.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
    NAGtestloss=np.mean(np.log(1+np.exp(-y_test.reshape(-1,1)*(X_test.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
    dw=lamda*(w-gamma*v) - np.mean((y_batch.reshape(-1,1) * X_batch)/(1 + np.exp(y_batch.reshape(-1,1) * X_batch.dot((w-gamma*v).T))), \
                          axis = 0, keepdims = True)
    v=gamma*v+learning_rate*dw
    w=w-v
    loss_NAGhistory.append(NAGloss)
    test_NAGhistory.append(NAGtestloss)

print(Accuracy(X_test,w))
```
0.8506848473680978

## RMSProp:

```
1   #  RMSProp
2
3   learning_rate=0.005
4   gamma=0.9
5   epsilon=1e-8
6   loss_RMSProphistory=[]
7   test_RMSProphistory=[]
8   w=0.001*np.random.random((1,124))
9   G=np.zeros((1,124))
10  for k in range(time):
11      mask=np.random.choice(X_train.shape[0],256,replace=False)
12      X_batch=X_train[mask]
13      y_batch=y_train[mask]
14      RMSProploss=np.mean(np.log(1+np.exp(-y_batch.reshape(-1,1)*(X_batch.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
15      RMSProptestloss=np.mean(np.log(1+np.exp(-y_test.reshape(-1,1)*(X_test.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
16
17      dw=lamda*w - np.mean((y_batch.reshape(-1,1) * X_batch)/(1 + np.exp(y_batch.reshape(-1,1) * X_batch.dot(w.T))), \
18                      axis = 0, keepdims = True)
19      G=gamma*G+(1-gamma)*np.multiply(dw,dw)
20
21      w=w-np.multiply(learning_rate/np.sqrt(G+epsilon),dw)
22
23      loss_RMSProphistory.append(RMSProploss)
24      test_RMSProphistory.append(RMSProptestloss)
25
26
27  print(Accuracy(X_test,w))
```
0.8503777409250046

## Adadelta:

```
1   #  AdaDelta
2
3   gamma=0.95
4   epsilon=1e-5
5   loss_AdaDeltahistory=[]
6   test_AdaDeltahistory=[]
7   w=0.001*np.random.random((1,124))
8   G=np.zeros((1,124))
9   delta_t=np.zeros((1,124))
10  for k in range(time):
11      mask=np.random.choice(X_train.shape[0],256,replace=False)
12      X_batch=X_train[mask]
13      y_batch=y_train[mask]
14      AdaDeltaloss=np.mean(np.log(1+np.exp(-y_batch.reshape(-1,1)*(X_batch.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
15      AdaDeltatestloss=np.mean(np.log(1+np.exp(-y_test.reshape(-1,1)*(X_test.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
16      dw=lamda*w - np.mean((y_batch.reshape(-1,1) * X_batch)/(1 + np.exp(y_batch.reshape(-1,1) * X_batch.dot(w.T))), \
17                      axis = 0, keepdims = True)
18      G=gamma*G+(1-gamma)*np.multiply(dw,dw)
19      delta_w=-np.multiply(np.sqrt(delta_t+epsilon)/np.sqrt(G+epsilon),dw)
20      w=w+delta_w
21      delta_t=gamma*delta_t+(1-gamma)*np.multiply(delta_w,delta_w)
22      loss_AdaDeltahistory.append(AdaDeltaloss)
23      test_AdaDeltahistory.append(AdaDeltatestloss)
24
25
26  print(Accuracy(X_test,w))
```
0.8517290092746146

Adam:

```
1   #  Adam
2
3   beta=0.9
4   gamma=0.999
5   eta=0.001
6   epsilon=1e-8
7   loss_Adamhistory=[]
8   test_Adamhistory=[]
9   w=0.001*np.random.random((1,124))
10  m=np.zeros((1,124))
11  G=np.zeros((1,124))
12  for k in range(time):
13      mask=np.random.choice(X_train.shape[0],4096,replace=False)
14      X_batch=X_train[mask]
15      y_batch=y_train[mask]
16      Adamloss=np.mean(np.log(1+np.exp(-y_batch.reshape(-1,1)*(X_batch.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
17      Adamtestloss=np.mean(np.log(1+np.exp(-y_test.reshape(-1,1)*(X_test.dot(w.T)))))+lamda/2*(np.linalg.norm(w, 2)**2)
18      dw=lamda*w - np.mean((y_batch.reshape(-1,1) * X_batch)/(1 + np.exp(y_batch.reshape(-1,1) * X_batch.dot(w.T))), \
19                           axis = 0, keepdims = True)
20      m=beta*m+(1-beta)*dw
21      G=gamma*G+(1-gamma)*np.multiply(delta_w,delta_w)
22      alpha=learning_rate*(np.sqrt(1-gamma**(k+1)))/(1-beta**(k+1))
23      w=w-alpha*m/(np.sqrt(G+epsilon))
24      loss_Adamhistory.append(Adamloss)
25      test_Adamhistory.append(Adamtestloss)
26
27
28  print(Accuracy(X_test,w))
```

0.8510533750998096

Linear Classification and Stochastic Gradient Descent:

NAG:

```
1   #  NAG
2
3   v=np.zeros((1,124))
4   gamma=0.9
5   w=np.random.random((1,124))*0.001
6   learning_rate=0.05
7   loss_NAGhistory=[]
8   test_NAGhistory=[]
9
10  for k in range(time):
11      mask=np.random.choice(X_train.shape[0],256,replace=False)
12      X_batch=X_train[mask]
13      y_batch=y_train[mask]
14      dw=np.zeros((1,X_batch.shape[1]))
15      NAGloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_batch.reshape(-1,1)*(X_batch.dot(w.T)),0))
16      NAGtestloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_test.reshape(-1,1)*(X_test.dot(w.T)),0))
17      loss_NAGhistory.append(NAGloss)
18      test_NAGhistory.append(NAGtestloss)
19      for i in range(X_batch.shape[0]):
20          if 1-y_batch[i]*(X_batch[i].dot((w-gamma*v).T))>=0:
21              dw += -y_batch[i]*X_batch[i].reshape((1,124))
22      dw= dw/X_batch.shape[0]+ lamda*(w-gamma*v)
23      v=gamma*v+learning_rate*dw
24      w=w-v
25
26  print(Accuracy(X_test,w))
```

0.8480437319574965

## RMSProp:

```python
# RMSProp

gamma=0.95
epsilon=1e-8
learning_rate=0.005
w=np.random.random((1,124))*0.001
G=np.zeros((1,124))
loss_RMSProphistory=[]
test_RMSProphistory=[]

for k in range(time):
    mask=np.random.choice(X_train.shape[0],256,replace=False)
    X_batch=X_train[mask]
    y_batch=y_train[mask]
    dw=np.zeros((1,X_batch.shape[1]))
    RMSProploss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_batch.reshape(-1,1)*(X_batch.dot(w.T)),0))
    RMSProptestloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_test.reshape(-1,1)*(X_test.dot(w.T)),0))
    loss_RMSProphistory.append(RMSProploss)
    test_RMSProphistory.append(RMSProptestloss)
    for i in range(X_batch.shape[0]):
        if 1-y_batch[i]*(X_batch[i].dot(w.T))>=0:
            dw += -y_batch[i]*X_batch[i].reshape((1,124))
    dw= dw/X_batch.shape[0]+ lamda*w
    G=gamma*G+(1-gamma)*np.multiply(dw,dw)
    w=w-np.multiply(learning_rate/np.sqrt(G+epsilon),dw)

print(Accuracy(X_test,w))
```

0.8502548983477674

## Adadelta:

```python
# AdaDelta

gamma=0.95
epsilon=1e-5
G=np.zeros((1,124))
delta_t=np.zeros((1,124))
loss_AdaDeltahistory=[]
test_AdaDeltahistory=[]
w=np.random.random((1,124))*0.001
for k in range(time):
    mask=np.random.choice(X_train.shape[0],512,replace=False)
    X_batch=X_train[mask]
    y_batch=y_train[mask]
    dw=np.zeros((1,X_batch.shape[1]))
    AdaDeltaloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_batch.reshape(-1,1)*(X_batch.dot(w.T)),0))
    AdaDeltatestloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_test.reshape(-1,1)*(X_test.dot(w.T)),0))
    loss_AdaDeltahistory.append(AdaDeltaloss)
    test_AdaDeltahistory.append(AdaDeltatestloss)
    for i in range(X_batch.shape[0]):
        if 1-y_batch[i]*(X_batch[i].dot(w.T))>=0:
            dw += -y_batch[i]*X_batch[i].reshape((1,124))
    dw= dw/X_batch.shape[0]+ lamda*w
    G=gamma*G+(1-gamma)*np.multiply(dw,dw)
    delta_w=-np.multiply(np.sqrt(delta_t+epsilon)/np.sqrt(G+epsilon),dw)
    w=w+delta_w
    delta_t=gamma*delta_t+(1-gamma)*np.multiply(delta_w,delta_w)

print(Accuracy(X_test,w))
```

0.8482279958233524

Adam:

```
1   # Adam
2
3   beta=0.9
4   gamma=0.999
5   eta=0.001
6   epsilon=1e-8
7   learning_rate=0.001
8   m=np.zeros((1,124))
9   G=np.zeros((1,124))
10  loss_Adamhistory=[]
11  test_Adamhistory=[]
12  w=np.random.random((1,124))*0.001
13  for k in range(time):
14      mask=np.random.choice(X_train.shape[0],4096,replace=False)
15      X_batch=X_train[mask]
16      y_batch=y_train[mask]
17      dw=np.zeros((1,X_batch.shape[1]))
18      Adamloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_batch.reshape(-1,1)*(X_batch.dot(w.T)),0))
19      Adamtestloss = lamda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_test.reshape(-1,1)*(X_test.dot(w.T)),0))
20      loss_Adamhistory.append(Adamloss)
21      test_Adamhistory.append(Adamtestloss)
22      for i in range(X_batch.shape[0]):
23          if 1-y_batch[i]*(X_batch[i].dot(w.T))>=0:
24              dw += -y_batch[i]*X_batch[i].reshape((1,124))
25      dw= dw/X_batch.shape[0]+ lamda*w
26      m=beta*m+(1-beta)*dw
27      G=gamma*G+(1-gamma)*np.multiply(delta_w,delta_w)
28      alpha=learning_rate*(np.sqrt(1-gamma**(k+1)))/(1-beta**(k+1))
29      w=w-alpha*m/(np.sqrt(G+epsilon))
30
31  print(Accuracy(X_test,w))
```
0.8494564215957251

## 8. The initialization method of model parameters:

W: set with random numbers and multiply 0.001;

G, v, m,　delta_t :set all to zeros.

## 9. The selected loss function and its derivatives:

Logistic Regression and Stochastic Gradient Descent：

$$\mathcal{L} = \frac{1}{n}\sum_{i=1}^{n} log(1 + e^{-y_i w^T x_i}) + \frac{\lambda}{2}\|w\|_2^2$$

$$\frac{\partial \mathcal{L}}{\partial w} = -\frac{1}{n}\sum_{i=1}^{n}\frac{y_i x_i}{1 + e^{y_i w^T x_i}} + \lambda w$$

Linear Classification:

$$\mathcal{L} = \frac{\lambda\|w\|^2}{2} + \frac{1}{n}\sum_{i=1}^{n} max(0, 1 - yi(w^T xi))$$

$$g_w(xi) = -yixi \quad 1 - yi(w^T xi) \geq 0$$

$$g_w(xi) = 0 \quad 1 - yi(w^T xi) < 0$$

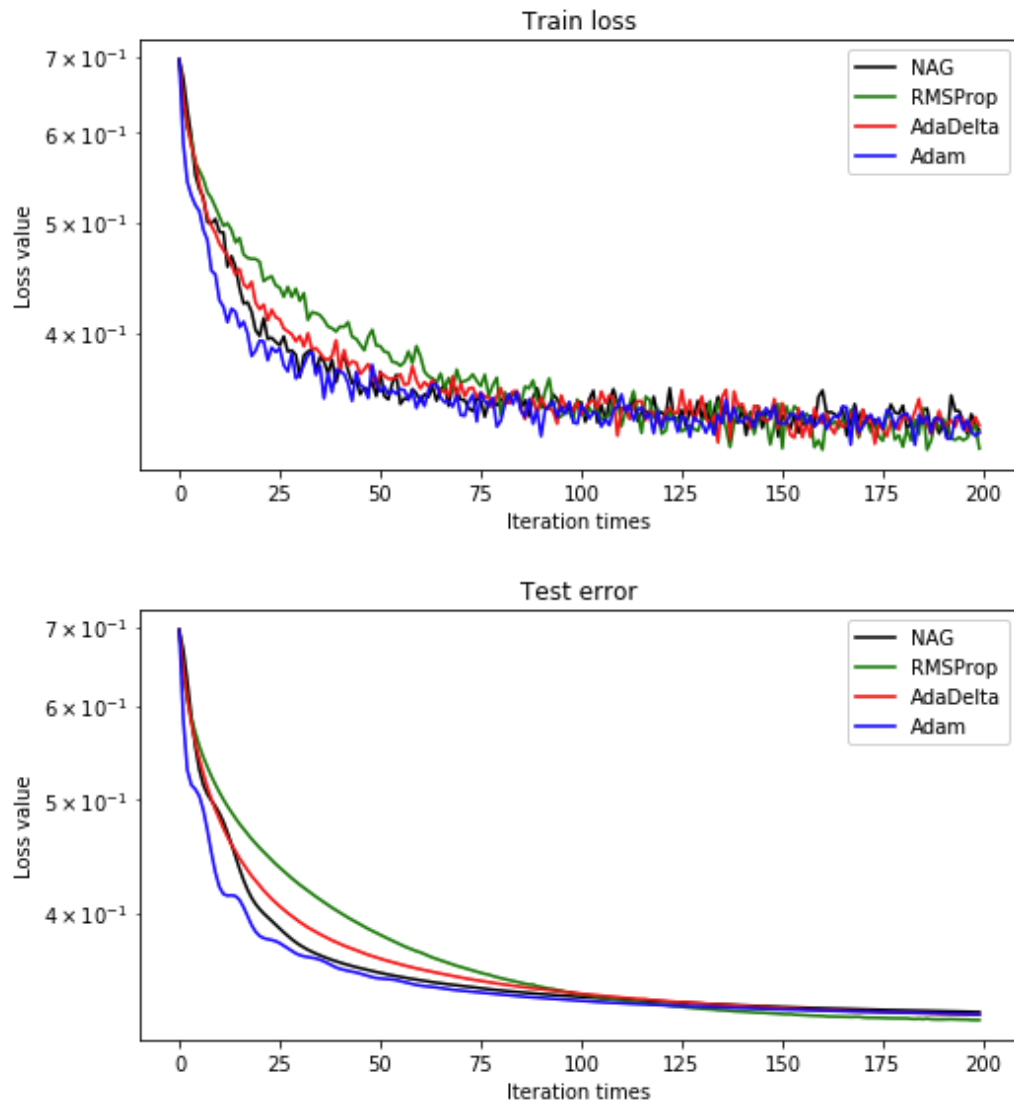$$\frac{\partial \mathcal{L}}{\partial w} = \frac{1}{n}\sum_{i=1}^{n} g_w(xi) + \lambda w$$

## 10. Experimental results and curve:

Logistic Regression and Stochastic Gradient Descent：

$\lambda = 0$, Iteration times=200, Batch_size=4096

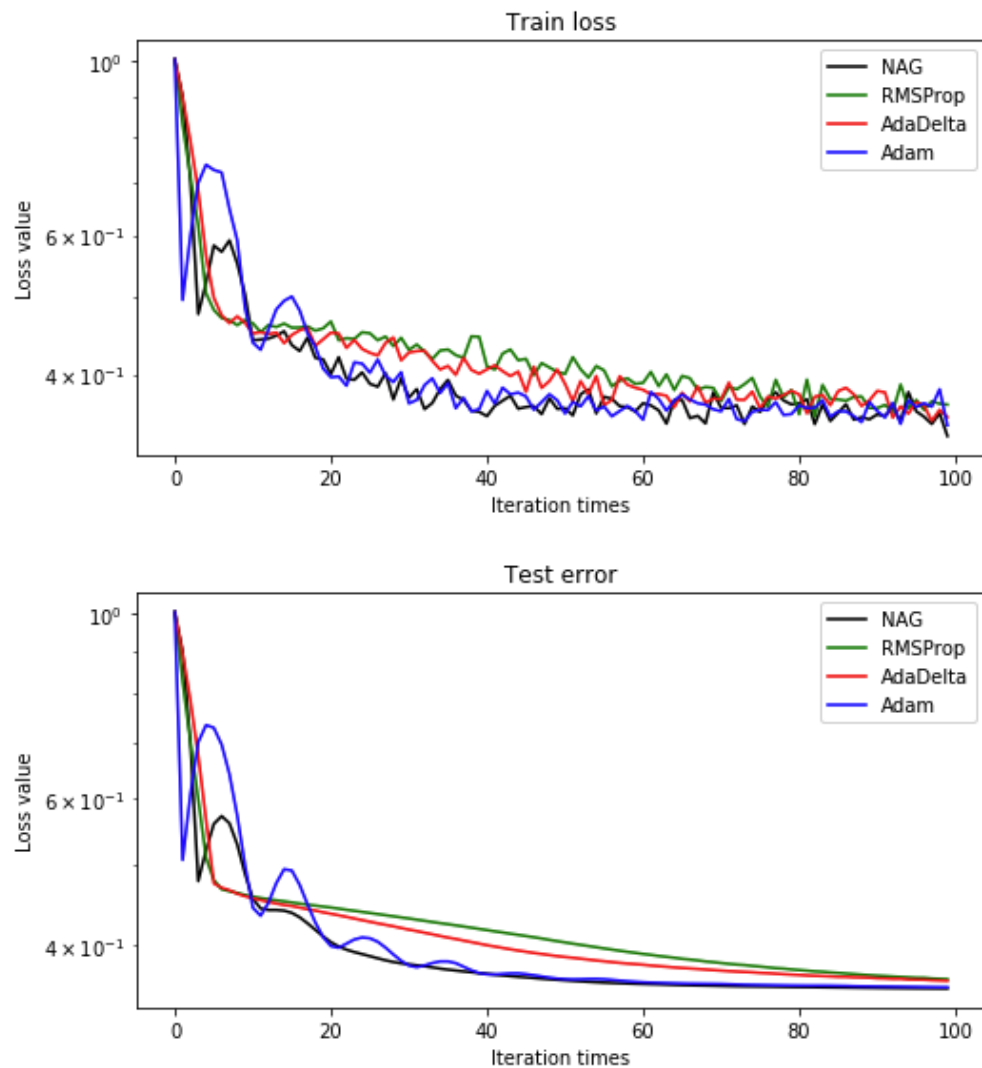|  | Hyper-parameter selection | Predicted Results (Best Results) |
|---|---|---|
| NAG | gamma=0.9<br>learning_rate=0.05 | 23 print(Accuracy(X_test,w))<br>0.8506848473680978 |
| RMSProp | learning_rate=0.005<br>gamma=0.9<br>epsilon=1e-8 | 27 print(Accuracy(X_test,w))<br>0.8503777409250046 |
| AdaDelta | gamma=0.95<br>epsilon=1e-5 | 26 print(Accuracy(X_test,w))<br>0.8517290092746146 |
| Adam | 3 beta=0.9<br>4 gamma=0.999<br>5 eta=0.001<br>6 epsilon=1e-8<br>7 learning_rate=0.001 | 28 print(Accuracy(X_test,w))<br>0.8510533750998096 |

Loss curve:





Linear Classification and Stochastic Gradient Descent:

$\lambda = 0$, Iteration times=100, Batch_size=4096

| | Hyper-parameter selection | Predicted Results (Best Results) |
|---|---|---|
| NAG | gamma=0.9<br>learning_rate=0.05 | 27 print(Accuracy(X_test, w))<br>0.8480437319574965 |
| RMSProp | 3 gamma=0.95<br>4 epsilon=1e-8<br>5 learning_rate=0.005 | 27 print(Accuracy(X_test, w))<br>0.8502548983477674 |
| AdaDelta | gamma=0.95<br>epsilon=1e-5 | 28 print(Accuracy(X_test, w))<br>0.8482279958233524 |

| Adam | 3  beta=0.9<br>4  gamma=0.999<br>5  eta=0.001<br>6  epsilon=1e-8<br>7  learning_rate=0.001 | 31  print(Accuracy(X_test,w))<br><br>0.8494564215957251 |
| --- | --- | --- |

Loss curve:



Train loss



Test error

## 11.Results analysis:

1. When the Regular modulus equal to 0, the accuracy is the highest. And the loss is much smaller.

2. Batch size is association with the amplitude of train loss, the larger the batch size, the smaller the amplitude of train loss.

3. If learning_rate in NAG =0.001, NAG should iterate a lot of time to convergence(probably 3000times), so I change the learning_rate to 0.05 in NAG. And the iteration time reduce a lot. So do the RMSProp, so I change the learning_rate to 0.005 in RMSProp.

4. If epsilon in AdaDelta = 1e-8, AdaDelta should iterate a lot of time to convergence(probably 3000times), so I change the epsilon to 1e-5, it perform much better.

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

Both methods are common classification algorithms. Look at the objective function, the difference is that logistic loss is used in logistic regression, svm uses hinge loss. The purpose of these two loss functions is to increase the weight of the data points that have a greater impact on the classification and reduce the weight of the data points with less classification. SVM approach is to consider only support vectors, which is the most relevant and classification of the few points to learn classifier. Logistic regression reduces the weights of points farther away from the classification plane through nonlinear mapping, and increases the weight of the data points most relevant to the classification.

## 13. Summary:

The algorithm behaves differently on different models and also has a relationship with hyperparameter learning rate $\eta$ and so on.