



华南理工大学

South China University of Technology

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:

Logistic Regression, Linear Classification and Stochastic Gradient Descent

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:

```
1 # NAG
2
3 learning_rate=0.05
4 w=0.001*np.random.random((1,124))
5 v=np.zeros((1,124))
6 gamma=0.9
7
8 loss_NAGhistory=[]
9 test_NAGhistory=[]
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     NAGloss=np.mean(np.log(1+np.exp(-y_batch.reshape(-1,1)*(X_batch.dot(w.T))))) + lambda/2*(np.linalg.norm(w, 2)**2)
15     NAGtestloss=np.mean(np.log(1+np.exp(-y_test.reshape(-1,1)*(X_test.dot(w.T))))) + lambda/2*(np.linalg.norm(w, 2)**2)
16     dw=lambda*(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))), \
17                                     axis = 0, keepdims = True)
18     v=gamma*v+learning_rate*dw
19     w=w-v
20     loss_NAGhistory.append(NAGloss)
21     test_NAGhistory.append(NAGtestloss)
22
23 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)))+lamba/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)))+lamba/2*(np.linalg.norm(w, 2)**2))
16
17     dw=lamba*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)))+lamba/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)))+lamba/2*(np.linalg.norm(w, 2)**2))
16     dw=lamba*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)))+landa/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)))+landa/2*(np.linalg.norm(w, 2)**2)
18     dw=landa*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 = landa * 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 = landa * 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]+ landa*(w-gamma*v)
23     v=gamma*v+learning_rate*dw
24     w=w-v
25
26 print(Accuracy(X_test,w))

```

0.8480437319574965

RMSProp:

```

1  # RMSProp
2
3  gamma=0.95
4  epsilon=1e-8
5  learning_rate=0.005
6  w=np.random.random((1,124))*0.001
7  G=np.zeros((1,124))
8  loss_RMSPropHistory=[]
9  test_RMSPropHistory=[]
10
11  for k in range(time):
12      mask=np.random.choice(X_train.shape[0],256,replace=False)
13      X_batch=X_train[mask]
14      y_batch=y_train[mask]
15      dw=np.zeros((1,X_batch.shape[1]))
16      RMSPropLoss = lambda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_batch.reshape(-1,1)*(X_batch.dot(w.T)),0))
17      RMSPropTestLoss = lambda * 0.5*(np.linalg.norm(w,2)**2)+np.mean(np.maximum(1-y_test.reshape(-1,1)*(X_test.dot(w.T)),0))
18      loss_RMSPropHistory.append(RMSPropLoss)
19      test_RMSPropHistory.append(RMSPropTestLoss)
20      for i in range(X_batch.shape[0]):
21          if 1-y_batch[i]*(X_batch[i].dot(w.T))>=0:
22              dw += -y_batch[i]*X_batch[i].reshape((1,124))
23      dw=dw/X_batch.shape[0]+ lambda*w
24      G=gamma*G+(1-gamma)*np.multiply(dw,dw)
25      w=w-np.multiply(learning_rate/np.sqrt(G+epsilon),dw)
26
27  print(Accuracy(X_test,w))

```

0.8502548983477674

Adadelta:

```

1  # AdaDelta
2
3  gamma=0.95
4  epsilon=1e-5
5  G=np.zeros((1,124))
6  delta_t=np.zeros((1,124))
7  loss_AdaDeltaHistory=[]
8  test_AdaDeltaHistory=[]
9  w=np.random.random((1,124))*0.001
10  for k in range(time):
11      mask=np.random.choice(X_train.shape[0],512,replace=False)
12      X_batch=X_train[mask]
13      y_batch=y_train[mask]
14      dw=np.zeros((1,X_batch.shape[1]))
15      AdaDeltaLoss = lambda * 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      AdaDeltaTestLoss = lambda * 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_AdaDeltaHistory.append(AdaDeltaLoss)
18      test_AdaDeltaHistory.append(AdaDeltaTestLoss)
19      for i in range(X_batch.shape[0]):
20          if 1-y_batch[i]*(X_batch[i].dot(w.T))>=0:
21              dw += -y_batch[i]*X_batch[i].reshape((1,124))
22      dw=dw/X_batch.shape[0]+ lambda*w
23      G=gamma*G+(1-gamma)*np.multiply(dw,dw)
24      delta_w=np.multiply(np.sqrt(delta_t+epsilon)/np.sqrt(G+epsilon),dw)
25      w=w+delta_w
26      delta_t=gamma*delta_t+(1-gamma)*np.multiply(delta_w,delta_w)
27
28  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 n=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*n+(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*n/(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:

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

$$g_w(x_i) = -y_i x_i \quad 1 - y_i(w^T x_i) \geq 0$$

$$g_w(x_i) = 0 \quad 1 - y_i(w^T x_i) < 0$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_w(\mathbf{x}_i) + \lambda \mathbf{w}$$

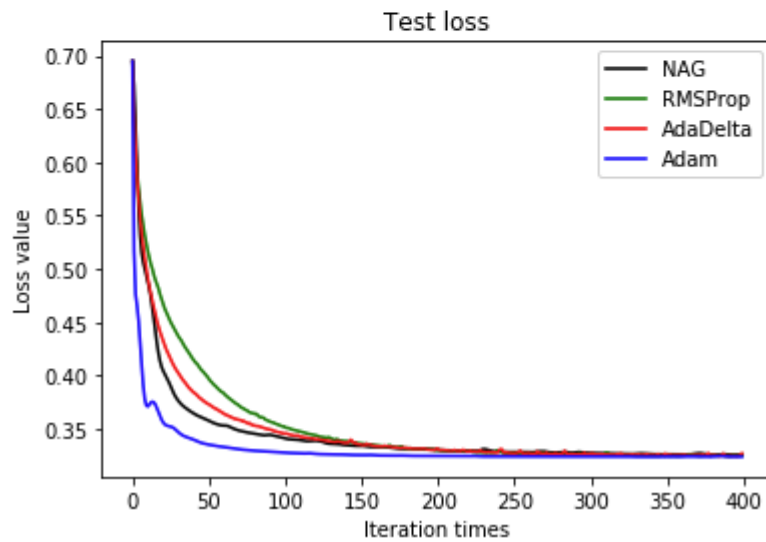
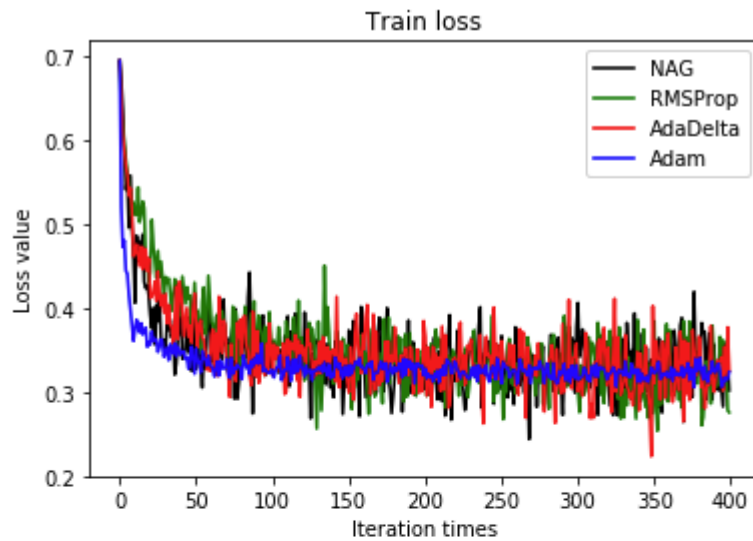
10. Experimental results and curve:

Logistic Regression and Stochastic Gradient Descent:

$\lambda = 0$, Iteration times=400

	Hyper-parameter selection	Predicted Results (Best Results)
NAG	<pre>gamma=0.9 learning_rate=0.05</pre> <p>Batch_size=256</p>	<pre>23 print(Accuracy(X_test,w))</pre> <p>0.8506848473680978</p>
RMSProp	<pre>learning_rate=0.005 gamma=0.9 epsilon=1e-8</pre> <p>Batch_size=256</p>	<pre>27 print(Accuracy(X_test,w))</pre> <p>0.8503777409250046</p>
AdaDelta	<pre>gamma=0.95 epsilon=1e-5</pre> <p>Batch_size=256</p>	<pre>26 print(Accuracy(X_test,w))</pre> <p>0.8517290092746146</p>
Adam	<pre>beta=0.9 gamma=0.999 eta=0.001 epsilon=1e-8</pre> <p>Batch_size=4096</p>	<pre>28 print(Accuracy(X_test,w))</pre> <p>0.8510533750998096</p>

Loss curve:



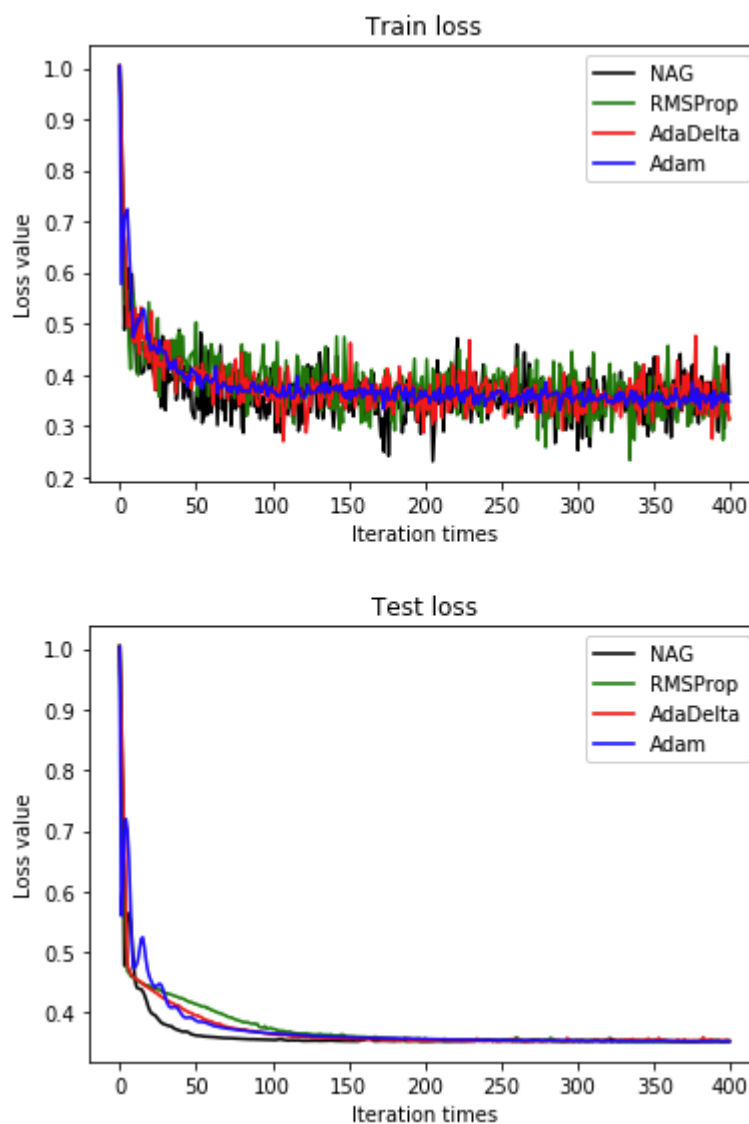
Linear Classification and Stochastic Gradient Descent:

$\lambda = 0$, Iteration times=400

	Hyper-parameter selection	Predicted Results (Best Results)
NAG	<pre>gamma=0.9 learning_rate=0.05</pre> <p>Batch_size=256</p>	<pre>27 print(Accuracy(X_test,w))</pre> <p>0.8480437319574965</p>
RMSProp	<pre>3 gamma=0.95 4 epsilon=1e-8 5 learning_rate=0.005</pre> <p>Batch_size=256</p>	<pre>27 print(Accuracy(X_test,w))</pre> <p>0.8502548983477674</p>

AdaDelta	<pre>gamma=0.95 epsilon=1e-5</pre> <p>Batch_size=512</p>	<pre>28 print(Accuracy(X_test,w))</pre> <p>0.8482279958233524</p>
Adam	<pre>3 beta=0.9 4 gamma=0.999 5 eta=0.001 6 epsilon=1e-8 7 learning_rate=0.001</pre> <p>Batch_size=4096</p>	<pre>31 print(Accuracy(X_test,w))</pre> <p>0.8494564215957251</p>

Loss curve:



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 3000 times), 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 3000 times), 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 η and so on.