

---

# CSE253 Homework 1

## Enroll date : 9pm 15 Jan.

---

**Zhenchao Gan, A53092819**  
Department of Computer Science  
9450 Gilman DR, La Jolla, CA  
zhgan@eng.ucsd.edu

### 1 Problems from Bishop

#### 1.1

$$\prod_{i=1}^d \int_{-\infty}^{\infty} e^{-x_i^2} dx_i = S_d \int_0^{\infty} e^{-r^2} r^{d-1} dr$$
$$S_d = \frac{\prod_{i=1}^d \int_{-\infty}^{\infty} e^{-x_i^2} dx_i}{\int_0^{\infty} e^{-r^2} r^{d-1} dr} \quad (1)$$

According to (1.41), we have

$$\int_{-\infty}^{\infty} \exp\{-\frac{2}{\lambda}x^2\}dx = (\frac{2\pi}{\lambda})^{1/2} \quad (2)$$

In (2), we set  $\lambda = 2$ , we can get

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \pi^{1/2}$$
$$\prod_{i=1}^d \int_{-\infty}^{\infty} e^{-x_i^2} dx_i = \pi^{d/2} \quad (3)$$

According to (1.44), we have

$$\Gamma(x) \equiv \int_0^{\infty} u^{x-1} e^{-u} du \quad (4)$$

In (4), we set  $u = r^2$ , we can get

$$\Gamma(x) \equiv \int_0^{\infty} r^{2(x-1)} e^{-r^2} dr^2 \equiv 2 \int_0^{\infty} r^{2x-1} e^{-r^2} dr \quad (5)$$

In (5), set  $d = 2x$

$$\Gamma(\frac{d}{2}) \equiv 2 \int_0^{\infty} r^{d-1} e^{-r^2} dr \quad (6)$$

Put (3) and (6) into (1), we can get

$$S_d = \frac{\prod_{i=1}^d \int_{-\infty}^{\infty} e^{-x_i^2} dx_i}{\int_0^{\infty} e^{-r^2} r^{d-1} dr}$$
$$= \frac{2\pi^{d/2}}{\Gamma(\frac{d}{2})} \quad (7)$$

When  $d = 2$ ,

$$\begin{aligned} S_2 &= \frac{2\pi^{d/2}}{\Gamma(\frac{d}{2})} \\ &= \frac{2\pi}{\Gamma(1)} \\ &= 2\pi \end{aligned} \quad (8)$$

The area of circle is  $\pi r^2$

When  $d = 3$ ,

$$\begin{aligned} S_3 &= \frac{2\pi^{d/2}}{\Gamma(\frac{d}{2})} \\ &= \frac{2\pi^{3/2}}{\Gamma(\frac{3}{2})} \\ &= 4\pi \end{aligned} \quad (9)$$

The area of sphere is  $\frac{4}{3}\pi r^3$

## 1.2

For a hypersphere with radius  $a$  and dimension  $d$ , the surface  $S = S_d a^{d-1}$ . Assuming a small cube in unit sphere with length of side  $L$ . Then the area is  $L^{d-1}$ . If the radius is  $a$ , then the length of new cube is  $aL$ , the new corresponding area is  $aL^{d-1}$ .

So

$$\begin{aligned} V_d &= \int_0^a S_d r^{d-1} dr \\ &= \frac{S_d}{d} \int_0^a r^d dr \\ &= \frac{S_d a^d}{d} \end{aligned} \quad (10)$$

$$\begin{aligned} \frac{\text{volume of Sphere}}{\text{volume of Cube}} &= \frac{2\pi^{\frac{d}{2}} a^d}{\Gamma(d/2) d (2a)^d} \\ &= \frac{\pi^{\frac{d}{2}}}{d 2^{d-1} \Gamma(d/2)} \end{aligned} \quad (11)$$

When  $d \rightarrow \infty$

$$\begin{aligned} \frac{\text{volume of Sphere}}{\text{volume of Cube}} &= \frac{\pi^{\frac{d}{2}}}{d 2^{d-1} \Gamma(d/2)} \\ &= \frac{\pi e^{d/2-1/2}}{d^{d/2-1/2}} \frac{e^{-1/2}}{d^{d/2-1/2}} (\text{substitute } \Gamma(d/2)) \\ &= 0 \cdot 0 (d \rightarrow \infty) \\ &= 0 \end{aligned} \quad (12)$$

Next, calculating ratio of the distance from the centre of the hypercube to one of the corners, divided by the perpendicular distance to one of the edges,

$$\begin{aligned} \frac{Dis(\text{corner})}{Dis(\text{edge})} &= \frac{\sqrt{\sum_{i=1}^d a^2}}{a} \\ &= \sqrt{d} \end{aligned} \quad (13)$$

It is straight-forward to see when  $d \rightarrow \infty$ , this value goes to  $\infty$

## 1.3

$$\begin{aligned} f &= 1 - \frac{V_{a-\epsilon}}{V_a} \\ &= 1 - \frac{(a-\epsilon)^d}{a^d} \\ &= 1 - \left(1 - \frac{\epsilon}{a}\right)^d \end{aligned} \quad (14)$$

When  $d \rightarrow \infty$

$$\begin{aligned} \lim_{d \rightarrow \infty} f &= \lim_{d \rightarrow \infty} 1 - (1 - \frac{\epsilon}{a})^d \\ &= 1 - 0 \\ &= 1 \end{aligned} \quad (15)$$

When  $\epsilon/a = 0.01$ ,

$$\begin{aligned} f_{d=2} &= 1 - 0.99^2 = 0.0199 \\ f_{d=10} &= 1 - 0.99^{10} = 0.0956 \\ f_{d=1000} &= 1 - 0.99^{1000} = 0.99996 \end{aligned}$$

When lies inside the radius  $a/2$ ,

$$\begin{aligned} f_{d=2} &= 0.5^2 = 0.25 \\ f_{d=10} &= 1 - 0.5^{10} = 0.000977 \\ f_{d=1000} &= 1 - 0.5^{1000} \approx 0.00000 \end{aligned}$$

#### 1.4

In Problem 1.2, we have show that for a hypersphere with radius  $a$  and dimension  $d$ , the surface  $S = S_d a^{d-1}$ .

So,

$$\begin{aligned} \int_{shell} p(x) dx &= p(r) S_d r^{d-1} \epsilon \\ &= \rho(r) \epsilon \end{aligned} \quad (16)$$

Next, we calculate the single maximum point,

$$\begin{aligned} \rho(r) &= \frac{S_d r^{d-1}}{(2\pi\sigma^2)^{1/2}} \exp(-\frac{r^2}{2\sigma^2}) \\ &\propto r^{d-1} \exp(-\frac{r^2}{2\sigma^2}) = f(r) \end{aligned} \quad (17)$$

We want to find the value of  $r$  to make  $f(r)$  reach maximum value, so

$$\begin{aligned} \frac{df(r)}{dr} &= [(d-1)r^{d-2} - r^{d-1} \frac{r}{\sigma^2}] \exp(-\frac{r^2}{2\sigma^2}) = 0 \\ \Rightarrow \hat{r} &= \sigma \sqrt{d-1} \\ \Rightarrow \hat{r} &\approx \sigma \sqrt{d} (d \gg 1) \end{aligned} \quad (18)$$

From equation (17), we have

$$\begin{aligned} \rho(r) &\propto r^{d-1} \exp(-\frac{r^2}{2\sigma^2}) \\ &= \exp((d-1) \ln r - \frac{r^2}{2\sigma^2}) \end{aligned} \quad (19)$$

So,

$$\begin{aligned} \frac{\rho(\hat{r}+\epsilon)}{\rho(\hat{r})} &= \exp^{(d-1) \ln(1+\frac{\epsilon}{\hat{r}}) - \frac{2\epsilon\hat{r}+\epsilon^2}{2\sigma^2}} \\ &= \exp^{(d-1)(\frac{\epsilon}{\hat{r}} - \frac{\epsilon^2}{2\hat{r}^2}) - \frac{2\epsilon\hat{r}+\epsilon^2}{2\sigma^2}} \\ &= \exp^{-\frac{\epsilon^2}{\sigma^2}} (substitute (d-1) with (\frac{\hat{r}}{\sigma})^2) \end{aligned} \quad (20)$$

## 2 Perceptron

### 2.1

$$\begin{aligned} y(x) &= 0 \\ \Rightarrow \omega^T x - \theta &= 0 \\ \Rightarrow \omega_1 x_1 + \omega_2 x_2 - \theta &= 0 \end{aligned} \quad (21)$$

The decision boundary line is  $\omega^T x = \theta$

$$\begin{aligned} l &= \frac{ax_0 + by_0 + c}{\sqrt{a^2 + b^2}} \\ &= \frac{-\theta}{\sqrt{a^2 + b^2}} \\ &= \frac{-\omega_0}{\|\omega\|} \\ &= \frac{-\omega_0}{\|\omega\|} \end{aligned} \quad (22)$$

### 2.2

(a)

Learning rule: If output is 1 and should be 0, then lower weights to active inputs and raise the threshold; If output is 0 and should be 1, then raise weights to active inputs and lower the threshold.

(b)

| $x_1$ | $x_2$ | Output | Teacher | $w_1$ | $w_2$ | Threshold ( $\theta$ ) |
|-------|-------|--------|---------|-------|-------|------------------------|
| 1     | 1     | 1      | 0       | -1    | -1    | 1                      |
| 0     | 1     | 0      | 1       | -1    | 0     | 0                      |
| 1     | 0     | 0      | 1       | 0     | 0     | -1                     |
| 1     | 1     | 1      | 0       | -1    | -1    | 0                      |
| 0     | 1     | 0      | 1       | -1    | 0     | -1                     |
| 1     | 1     | 1      | 0       | -2    | -1    | 0                      |
| 1     | 0     | 0      | 1       | -1    | -1    | -1                     |

$$w_1 = -1, w_2 = -1, \theta = -1$$

(c)

It is not unique.

For example,  $w_1 = -2, w_2 = -1, \theta = -2$

### 2.3

(i.)

The implementation is in the appendix.

The advantage of the z score transformation is that it takes into account both the mean value and the variability in a set of raw data.

(ii.)

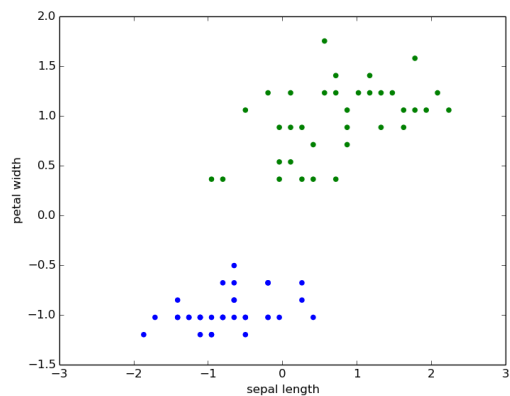
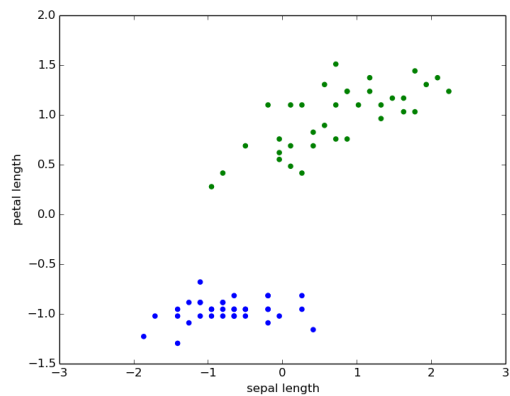
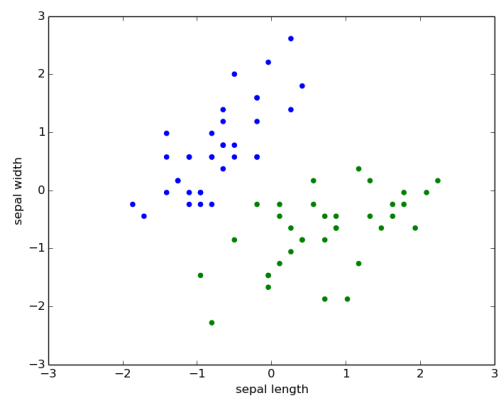
Blue colour stand for setosa, red colour stands for versicolor.

Classes linearly separable in each of the feature spaces. Because it is easy to find a boundary line in each figure below.

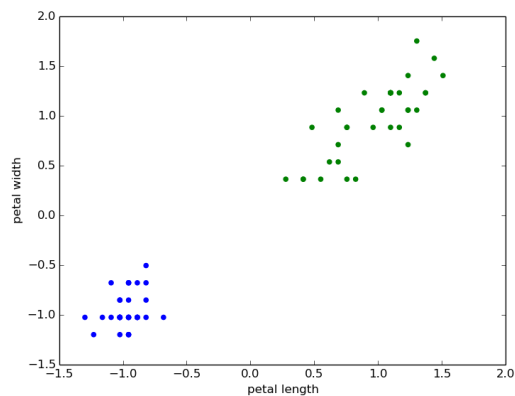
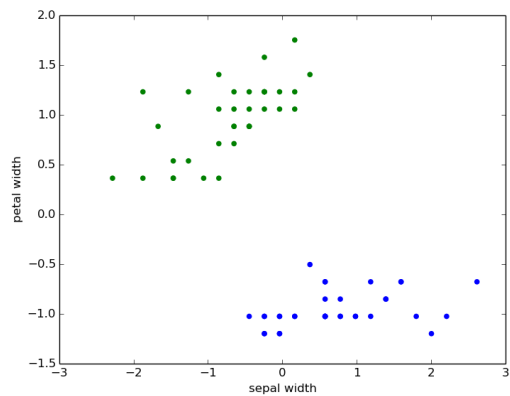
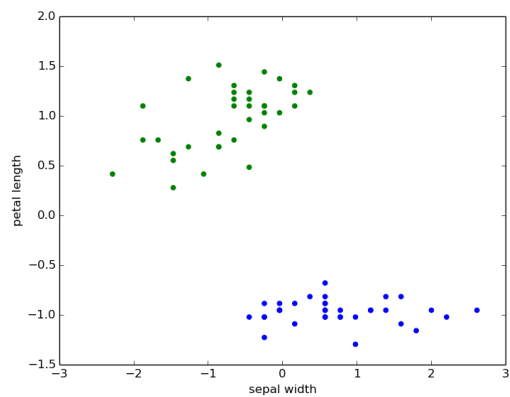
(iii.)

My source code is in the appendix.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269



270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323



(iv. v.)

| learning_rate | error rate |
|---------------|------------|
| .25           | 0.0%       |
| .5            | 0.0%       |
| 1.0           | 0.0%       |
| 2.0           | 0.0%       |

For this test set, our model and prediction is very perfect. I run many times and each iteration always choose a random point and try different learning rate, all get 0.0% error rate. So I think the data must be linearly separable. That is to say, the test data is very good.

### 3 Logistic and Softmax Regression

#### 3.1

$$\begin{aligned}
 -\frac{\partial E^n(\omega)}{\partial \omega_j} &= \frac{t^n}{g_w(x^n)} \frac{\partial g_w(x^n)}{\partial \omega_j} - \frac{(1-t^n)}{1-g_w(x^n)} \frac{\partial g_w(x^n)}{\partial \omega_j} \\
 &= t^n g_w(-x^n) x_j^n - (1-t^n) g_w(x^n) x_j^n \\
 &= x_j (-t^n (1-y^n) + (1-t^n) y^n) \\
 &= (t^n - y^n) x_j
 \end{aligned} \tag{23}$$

#### 3.2

First,  $\ln(\frac{a}{b}) = \ln a - \ln b \Rightarrow \ln y_k^n = a_k^n - \ln \sum_{k'} \exp(a_{k'}^n)$

Then,

$$\begin{aligned}
 -\frac{\partial E^n(\omega)}{\partial \omega_{jk}} &= \sum_{m=1}^c t_c^n \frac{\partial \ln y_m^n}{\partial \omega_{jk}} \\
 &= \sum_{m=1}^c t_c^n \frac{\partial a_m^n - \ln \sum_{k'} \exp(a_{k'}^n)}{\partial \omega_{jk}} \\
 &= t_k^n (x_j^n - \frac{\exp(a_k^n) x_j^n}{\sum_{k'} \exp(a_{k'}^n)}) + \sum_{m \neq k} t_m^n (-\frac{\exp(a_k^n) x_j^n}{\sum_{k'} \exp(a_{k'}^n)}) \\
 &= t_k^n x_j^n - \frac{\exp(a_k^n) x_j^n}{\sum_{k'} \exp(a_{k'}^n)} \sum_{m=1}^c t_m^n \\
 &= t_k^n x_j^n - \frac{\exp(a_k^n) x_j^n}{\sum_{k'} \exp(a_{k'}^n)} \\
 &= (t_k^n - y_k^n) x_j^n
 \end{aligned} \tag{24}$$

#### 3.3

The implementation is in appendix reference <http://g.sweyla.com/blog/2012/mnist-numpy/>.

#### 3.4

(a)

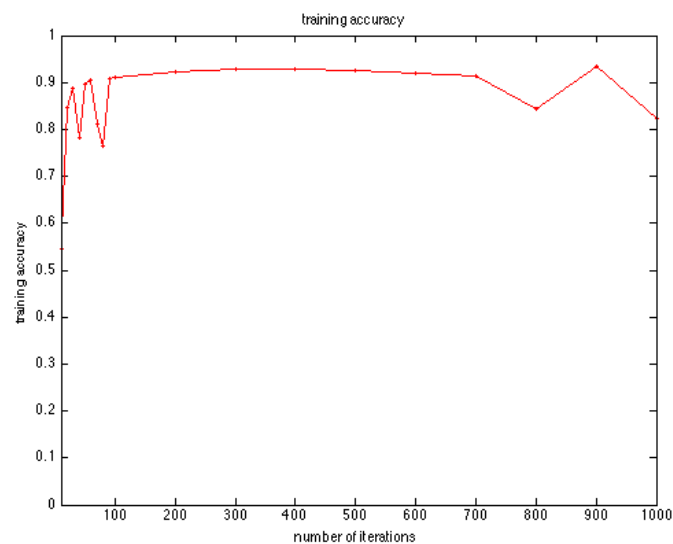
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

| Number | Accuracy |
|--------|----------|
| 0      | 98.3%    |
| 1      | 98.95%   |
| 2      | 97.3%    |
| 3      | 97.05%   |
| 4      | 97.4%    |
| 5      | 92.8%    |
| 6      | 97.8%    |
| 7      | 96.4%    |
| 8      | 91.9%    |
| 9      | 94.2%    |

(b)  
The overall accuracy is 82.25 %

## 4 Softmax Regression via Gradient Descent

(a)



(b)

The test accuracy on the test set is : 87.9%

(c)

The test accuracy is higher than the one-vs-all logistic regression approach.

## 5 Appendix

### 5.1 Perceptron Source Code

```
import numpy
import random
import matplotlib.pyplot as plt

train_data = []
```



```

432 label = ["sepal_length", "sepal_width", "petal_length", "petal_width"]
433 means = []
434 stds = []
435
436 class trainModel:
437
438     def __init__(self, learning_rate, dimension):
439         self.learning_rate = learning_rate
440         self.dimension = dimension
441         self.threshold = 0
442         self.weights = numpy.array([0 for i in range(dimension)])
443
444     def train(self, train_data):
445         self.ERROR = 0
446         Itertimes = 0
447         while Itertimes < 1000: #and not self.check(train_data):
448             row = train_data[random.randint(0, len(train_data)-1)]
449             teacher = row[self.dimension]
450             x = row[:self.dimension]
451             output = self.predict(x)
452             self.weights = self.weights + self.learning_rate * (teacher - output) * x
453             self.threshold = self.threshold - self.learning_rate * (teacher - output)
454             Itertimes = Itertimes + 1
455
456     def predict(self, x):
457         v = sum(x[:self.dimension] * self.weights)
458
459         if v >= self.threshold:
460             return 1
461         else:
462             return 0
463
464     def check(self, data):
465         err = 0
466         for row in data:
467             v = self.predict(row)
468             if v != row[self.dimension]:
469                 err = err + 1
470         if err == self.ERROR and err < 3:
471             return True
472         else:
473             self.ERROR = err
474             return False
475
476     def read(filepath):
477         file = open(filepath, "r")
478         data = []
479         while 1:
480             line = file.readline()
481             if not line:
482                 break
483             line = line.strip('\n')
484             row = line.split(',')
485             for i in xrange(4):
486                 row[i] = float(row[i])
487             row[4] = 1 if row[4] == "Iris-setosa" else 0
488             data.append(row)
489         return data

```

```

486 def zscore(data):
487     mydata = numpy.vstack(data)
488     for i in xrange(4):
489         mean = numpy.mean(mydata.T[i])
490         std = numpy.std(mydata.T[i])
491         means.append(mean)
492         stds.append(std)
493         for j in xrange(len(mydata.T[i])):
494             mydata[j][i] = (mydata[j][i] - mean)/std
495     return mydata
496
497 def plot_data():
498     setosa_data = []
499     versicolor_data = []
500     for ele in train_data:
501         if ele[4] == 1:
502             setosa_data.append(ele)
503         else:
504             versicolor_data.append(ele)
505     setosa_data = numpy.vstack(setosa_data)
506     versicolor_data = numpy.vstack(versicolor_data)
507
508     # plot my data
509     for i in xrange(3):
510         for j in xrange(3 - i):
511             plt.scatter(setosa_data.T[i], setosa_data.T[i+j+1], color = 'blue')
512             plt.scatter(versicolor_data.T[i], versicolor_data.T[i+j+1], color = 'green')
513             plt.xlabel(label[i])
514             plt.ylabel(label[i+j+1])
515             plt.show()
516
517 def Runtest(learning_rate):
518     print "Running as Learning Rate: ", learning_rate
519     model = trainModel(learning_rate, 4)
520     model.train(train_data)
521     data = read("iris_test.data")
522     mydata = numpy.vstack(data)
523     for i in xrange(4):
524         for j in xrange(len(mydata.T[i])):
525             mydata[j][i] = (mydata[j][i] - means[i])/stds[i]
526
527     miss = 0
528     for ele in mydata:
529         label = model.predict(ele)
530         if label != ele[-1]:
531             miss = miss+1
532
533     print "ERROR rate: ", miss * 100.0 / len(mydata), "%"
534
535 train_data = read("iris_train.data")
536 train_data = zscore(train_data)
537
538 Runtest(2)
539 Runtest(1)
540 Runtest(.5)
541 Runtest(.25)

```

## 5.2 Logistic and Softmax Regression Source Code

```

import os, struct
import numpy as np
from array import array as pyarray
from numpy import append, array, int8, uint8, zeros

def load_mnist(dataset="training", num = 20000, digits=np.arange(10), path="."):
    """
    Loads MNIST files into 3D numpy arrays

    Adapted from: http://abel.ee.ucla.edu/cvxopt/\_downloads/mnist.py
    """

    if dataset == "training":
        fname_img = os.path.join(path, 'train-images.idx3-ubyte')
        fname_lbl = os.path.join(path, 'train-labels.idx1-ubyte')
    elif dataset == "testing":
        fname_img = os.path.join(path, 't10k-images.idx3-ubyte')
        fname_lbl = os.path.join(path, 't10k-labels.idx1-ubyte')
    else:
        raise ValueError("dataset must be 'testing' or 'training'")

    flbl = open(fname_lbl, 'rb')
    magic_nr, size = struct.unpack(">II", flbl.read(8))
    lbl = pyarray("b", flbl.read())
    flbl.close()

    fimg = open(fname_img, 'rb')
    magic_nr, size, rows, cols = struct.unpack(">IIII", fimg.read(16))
    img = pyarray("B", fimg.read())
    fimg.close()

    ind = [ k for k in range(size) if lbl[k] in digits ]
    N = num

    images = zeros((N, rows*cols+1), dtype=uint8)
    labels = zeros(N, dtype=int8)
    for i in range(num):
        feature = img[ind[i]*rows*cols : (ind[i]+1)*rows*cols ]
        feature.insert(0, 1)
        images[i] = array(feature)
        labels[i] = lbl[ind[i]]

    return images, labels

class Logistic:

    def __init__(self, step, numlabel, dimension):
        self.step = step
        self.numlabel = numlabel
        self.weights = np.zeros(dimension)

    def sigmoid(self, data):
        return 1.0 / (1 + np.exp(-1 * (self.weights.dot(data))))

```

```

594     def train(self, traindata, label):
595         Itertime = 0
596         teacher = np.zeros(len(label))
597         for i in range(len(label)):
598             if label[i] == self.numlabel:
599                 teacher[i] = 1
600
601         while Itertime < 2000:
602             output = 1.0 / (1 + np.exp(-1.0 * (self.weights.dot(traindata.T) ) ) )
603             self.weights = self.weights + self.step * ( (teacher - output).dot(traindata.T) )
604             Itertime = Itertime + 1
605
606     def predict(self, x):
607         return self.sigmoid(x)
608
609 class softmax:
610
611     def __init__(self, step, dimension, nclass, data, labels):
612         self.step = step
613         self.dimension = dimension
614         self.nclass = nclass
615         self.data = data
616         self.labels = labels
617
618         self.weights = np.zeros((nclass, dimension))
619         self.teacher = np.zeros((nclass, len(labels)))
620         for i in range(len(labels)):
621             self.teacher[labels[i]][i] = 1
622
623     def train(self, Round):
624         Itertime = 0
625         while Itertime < Round:
626             print Itertime
627             output = self.predict(self.data)
628             self.weights = self.weights + self.step * (self.teacher - output).dot(self.data.T)
629             Itertime = Itertime + 1
630
631     def predict(self, data):
632         numerator = np.exp(self.weights.dot(data.T))
633         denominator = np.sum(numerator, axis = 0)
634         return numerator/denominator
635
636 def Lrun1():
637     for i in range(10):
638         model = Logistic(10e-8, i, 785)
639         print "Classify number", i
640         model.train(traindata, trainlabels)
641         N = 0
642         for index in range(len(testdata)):
643             v = model.predict(testdata[index])
644             if v >= 0.5 and testlabels[index] == i:
645                 N += 1
646             if v < 0.5 and testlabels[index] != i:
647                 N += 1
648         print "Classify number", i, "accuracy is: ", (N/2000.0)*100, "%"

```

```

648 def Lrun2():
649     models = []
650     for i in range(10):
651         model = Logistic(10e-9, i, 785)
652         model.train(traindata, trainlabels)
653         print "trainging", i, "finished"
654         models.append(model)
655         N = 0
656
657     for index in range(len(testdata)):
658         data = testdata[index]
659         m = -1.0
660         ind = -1
661         for j in range(10):
662             v = models[j].predict(data)
663             if v > m:
664                 m = v
665                 ind = j
666             if testlabels[index] == ind:
667                 N = N + 1
668         print "Overall Accuracy is : ", (N/2000.0)*100, "%"
669
670 def Srun1():
671     for i in range(2,11):
672         model = softmax(10e-9, 785, 10, traindata, trainlabels)
673         model.train(100*i)
674         output = model.predict(traindata)
675         output = np.argmax(output, axis = 0)
676         N = 0
677         for j in range(len(trainlabels)):
678             if output[j] == trainlabels[j]:
679                 N = N + 1
680         print "Iteration : ", i*100, " Accuracy is : ", (N/20000.0)*100, "%"
681
682 def Srun2():
683
684     model = softmax(10e-9, 785, 10, traindata, trainlabels)
685     model.train(500)
686     output = model.predict(testdata)
687     output = np.argmax(output, axis = 0)
688     N = 0
689     for j in range(len(testlabels)):
690         if output[j] == testlabels[j]:
691             N = N + 1
692     print " Accuracy is : ", (N*1.0/len(testlabels))*100, "%"
693
694 # read data
695 traindata, trainlabels = load_mnist('training', 20000)
696 testdata, testlabels = load_mnist('testing', 2000)
697
698 # run logistic
699 # Lrun1()
700 # Lrun2()
701
702 # run softmax
703 # Srun1()

```

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

|           |
|-----------|
| Srun2 ( ) |
|-----------|