Lab 7

# MACHINE LEARNING USING THE SCIKIT-LEARN AND THE MNIST DATASET

Xinyi(Blanca) Bai

Cornell University Electrical and Computer Engineering

# Table of Contents

# Part A Data Set Creation

Create 60 handwritten characters (50 of these are the digits 0-9 handwritten (5 of each) number; 10 of these to be symbols that are not the characters 0-9.

## 2. Photograph these characters into images.

I scan the whole big image, and then use screenshot to divide every character and number into a single image.

## 3. Preprocess the outline above.

First, threshold to make binary (bilevel) images and then resample to the required size and reposition according to the center of mass (COM).

The commands vfmt and vxport may be used to cornvert between VisionX image file format and other image file formats. Your final images should be stored in .png format Upload these images to Visionx, (use vfmt to import to VX format) and then:

```
1.  which vpix
2.  cd home/vision/v4
3.  ls
4.  cd templates
```

### 3.1. Convert them to gray-level images.
I scan the image, so it is already in gray-level, we can also use

### 3.2. Threshold the images to make them containing only background & foreground (binary images).

```
1.  /****************************************************************/
2.  /* threthold    Compute local max operation on a single byte image    */
3.  /****************************************************************/
4.
5.  #include "VisXV4.h"              /* VisionX structure include file    */
6.  #include "Vutil.h"              /* VisionX utility header files      */
7.
8.  VXparam_t par[] =              /* command line structure          */
9.  { /* prefix, value,   description                        */
10. {    "if=",    0,   " input file  vtemp: local max filter "},
11. {    "of=",    0,   " output file "},
12. {     0,      0,   0}  /* list termination */
13. };
14. #define  IVAL   par[0].val
15. #define  OVAL   par[1].val
16.
17. main(argc, argv)
18. int argc;
19. char *argv[];
20. {
21. Vfstruct (im);                      /* i/o image structure          */
22. Vfstruct (tm);                      /* temp image structure         */
23. int        y,x;                     /* index counters               */
24.   VXparse(&argc, &argv, par);       /* parse the command line       */
25.
26.   Vfread(&im, IVAL);                /* read image file              */
27.   Vfembed(&tm, &im, 1,1,1,1);       /* image structure with border  */
28.   if ( im.type != VX_PBYTE ) {      /* check image format           */
29.      fprintf(stderr, "vtemp: no byte image data in input file\n");
```

```
30.      exit(-1);
31.    }
32.    for (y = im.ylo ; y <= im.yhi ; y++) {   /* compute the function */
33.       for (x = im.xlo; x <= im.xhi; x++)  {      /***********************/
34.            //im.u[y][x] = MAX(tm.u[y][x],           /* You only need to     */
35.            //             MAX(tm.u[y][x+1],       /* change this section  */
36.            //               MAX(tm.u[y+1][x],    /* for your program     */
37.            //                 MAX(tm.u[y][x-1], /************************/
38.            //                   tm.u[y-1][x])))); 
39.
40.
41.            if(tm.u[y][x] >= 200 && tm.u[y][x] <= 255){
42.               im.u[y][x] = 255;
43.            }
44.
45.            else {
46.               im.u[y][x] = 0;
47.            }
48.        }
49.     }
50.    Vfwrite(&im, OVAL);                    /* write image file               */
51.    exit(0);
52. }
```

Compile and run my program

```
1.  vcc thre.c -o thre
2.  thre in=01.vx of=01th.vx
```

then new a file, write script to run my thre program on a batch of images(all images)

```
1.  loopthe
```

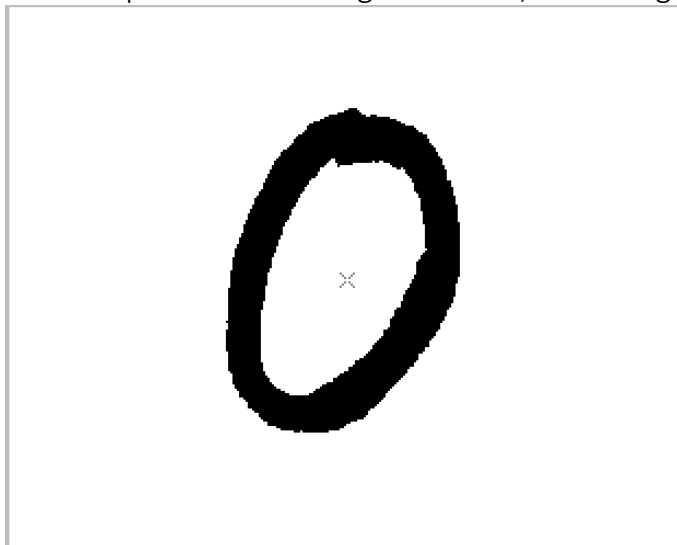I find I have no permission to run the script file
So I use chmod command to open all authority to that file

```
1.  chmod 744 loopthe
```

I set the pixel value of forground as 0, the background as 255.

## 3.3. For each image, generate a 20x20 bounding box that contains only the digit

one way to do this is to first do an exhaustive search of all the foreground pixel's coordinate, and bound all of the foreground pixels with a square box; then, you resize the square box into 20x20.

I wrote a program vcom.c to calculate the central mass of a image, and use the central mass value to crop images as 400X400 size.

vcom.c file

```
1.  /********************************************************************/
2.  /* vcom return the left lower corner of a byte image          */
3.  /********************************************************************/
4.  #include "VisXV4.h"
5.  #include "Vutil.h"
6.  VXparam_t par[] =
7.  { /* prefix, value,*/
8.  { "if=", 0, " input file vcorner: copy lower left corner"},
9.  { "of=", 0, " output file "},
10. { 0, 0, 0} /* list termination */
11. };
12.
13. #define IVAL par[0].val
14. #define OVAL par[1].val
15.
16. int main(argc, argv)
17. int argc;
18. char *argv[];
19. {
20. /* VisionX structure include file     */
21. /* VisionX utility header files       */
22. /* command line structure             */
23. /* description                        */
24.     Vfstruct (im);
25.     Vfstruct (tm);
26.     int y,x;
27.     int s;
28.     float bbx[6];
29.     int m00=0;
30.     int m01=0;
31.     int m10=0;
32.     int xx=0;
33.     int yy=0;
34.     VXparse(&argc, &argv, par); /* parse the command line */
35.     /* create VX image for result */
36.     s = 400; /* set size of result image */
37.     bbx[0] = bbx[2] = bbx[4] = bbx[5] = 0.0; bbx[1] = bbx[3] = s;
38.     Vfnewim(&tm, VX_PBYTE, bbx, 1);
39.     while ( Vfread(&im, IVAL) ) { /* read image file */
40.         if ( im.type != VX_PBYTE ) { /* check image format */
41.             fprintf(stderr, "vcorner: no byte image data in input file\n");
42.             exit(-1);
43.         }
44.
45.      /* check that the input image is large enough */
46.         if ( (im.xhi - im.xlo ) < (tm.xhi - tm.xlo) ||(im.yhi - im.ylo ) < (tm.yhi
    - tm.ylo) ) {
47.                 fprintf(stderr, "vcorner: input image too small\n");
48.             exit(-1);
49.         }
```

```
50.
51.
52.
53.          for(y = im.ylo; y <= im.yhi; y++){
54.              for(x = im.xlo; x <= im.xhi; x++){
55.              if( im.u[y][x] == 0){
56.                  m00++;
57.              }
58.              }
59.          }
60.
61.          for(y = im.ylo; y <= im.yhi; y++){
62.              for(x = im.xlo; x <= im.xhi; x++){
63.              if( im.u[y][x] == 0){
64.                  m01 = x + m01;
65.                  m10 = y + m10;
66.              }
67.              }
68.          }
69.
70.          xx = m01 / m00;
71.          yy = m10 / m00;
72.
73.          printf("xx = %d \n", xx);
74.          printf("yy = %d \n", yy);
75.          im.xlo = xx - 200;
76.          im.ylo = yy - 200;
77.
78.          for(y = tm.ylo; y <= tm.yhi; y++){
79.              for(x = tm.xlo; x <= tm.xhi; x++){
80.              tm.u[y][x] = im.u[y + im.ylo][x + im.xlo];
81.              }
82.          }
83.
84.
85.          Vfwrite(&tm, OVAL);
86.      }
87.          exit(0);
88. }
89. /* write image file                    */
```

Note： for imgae 54th.vx
The size of result image = 360
So I changed part of the code as below:

```
1.  Im.xlo = xx - 130;
2.  Im.ylo = yy - 130;
3.  Vmag if=54com.vx of=54scale.vx m=0.05556
```
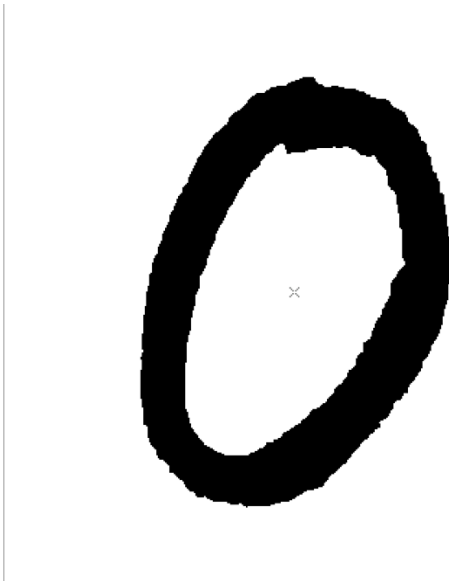
Compile and run my program

```
1.  vcc vcom.c -o vcom
2.  vcom if=02th.vx of=02com.vx
```

then new a file, write script to run my vcom program on a batch of images(all images)

```
1.  cp loopthe loopcom
```



Next, scale them to 20X20 size

```
1.  vmag if=02com.vx of=02scale.vx m=0.05
```

then new a file, write script to run vmag program on a batch of images(all images)

```
1.  cp loopthe loopmag
```



Next, reverse the color of image.

Note: because we need to add boundary to the 20X20 images, the color of the boundary should be the same as the background. My program to add boundary, the boundary I add is black color, so I do the color reverse at this step.

```
1.  Vpix -neg if=02scale.vx of=02inv.vx
2.  cp loopthe loopinv
```



## 3.4. Embed the boundary of that 20x20 image with background pixels to make it 28x28.

Write a program to add the boundary of the color-revered image.

```
1.  cp thre.c vboundary.c
2.  /*******************************************************    */
3.  /********************************************************/
4.
5.  #include "VisXV4.h"            /* VisionX structure include file    */
6.  #include "Vutil.h"             /* VisionX utility header files      */
7.
8.  VXparam_t par[] =             /* command line structure            */
9.  { /* prefix, value,    description                        */
10. {    "if=",   0,   " input file  vtemp: local max filter "},
11. {    "of=",   0,   " output file "},
12. {     0,      0,    0}  /* list termination */
13. };
14. #define  IVAL    par[0].val
15. #define  OVAL    par[1].val
16.
17. main(argc, argv)
18. int argc;
19. char *argv[];
20. {
21. Vfstruct (im);                          /* i/o image structure           */
22. Vfstruct (tm);                          /* temp image structure          */
23. int       y,x;                          /* index counters                */
24.   VXparse(&argc, &argv, par);           /* parse the command line        */
25.
26.   Vfread(&im, IVAL);                     /* read image file               */
27.   Vfembed(&tm, &im, 4,4,4,4);            /* image structure with border   */
28.   if ( im.type != VX_PBYTE ) {          /* check image format            */
29.      fprintf(stderr, "vtemp: no byte image data in input file\n");
30.      exit(-1);
31.   }
```

```
32.
33.
34.    Vfwrite(&tm, OVAL);              /* write image file              */
35.    exit(0);
36. }
```

Compile and run my program, using script to add boundary on all images.

```
1.  vcc vboundary.c -o vboundary
2.  vboundary if=02inv.vx of=02b.vx
3.  cp loopinc loopb
```



### 3.5  Export from VX to .png format using vxport.

Using script to convert format for all images.

```
1.  vxport if=02b.vx of=02b -png
2.  cp loopb looppng
```

## 4.Load the pre-processed data to scikit-learn using the python code we provided to you. By now you have the self-made train/test data for the classifier

The reason why we want you to first bound the digit in a 20x20 box before making it 28x28 is that we want to make sure that every digit you created will have the same size in the image (occupying 400/784 pixels). Assume that some of you may have very small-size handwriting, while some may

write the digit really large; in that way, even if you both shift the COM to the center, the size difference of your digit will still cause trouble for the classifier to learn.

```
1.  #1 import required modules
2.  import time
3.  import io
4.  import matplotlib.pyplot as plt
5.  import numpy as np
6.  from scipy.io.arff import loadarff
7.
8.  from sklearn.datasets import get_data_home
9.  from sklearn.externals.joblib import Memory
10. from sklearn.linear_model import LogisticRegression
11. from sklearn.model_selection import train_test_split
12. from sklearn.preprocessing import StandardScaler
13. from sklearn.utils import check_random_state
14. from urllib.request import urlopen
```

```
15. In [2]:
16. #2.  read the NMIST dataset
17. memory = Memory(get_data_home())
18. @memory.cache()
19. def fetch_mnist():
20.     content = urlopen(
21.         'https://www.openml.org/data/download/52667/mnist_784.arff').read()
22.     data, meta = loadarff(io.StringIO(content.decode('utf8')))
23.     data = data.view([('pixels', '<f8', 784), ('class', '|S1')])
24.     return data['pixels'], data['class']
25. X, y = fetch_mnist()
```

```
26. In [3]:
27. # rescale the data, use the traditional train/test split
28.
29. X = X / 255.
30.
31. ###### NEW  Refromat the the labels to be integers rather than byte arrays
32. y_trans = []
33. for i in range(len(y)):
34.     y_trans.append(int(y[i]))
35. y = np.asarray(y_trans)
36.
37. X_train, X_test = X[:60000], X[60000:]
38. y_train, y_test = y[:60000], y[60000:]
```
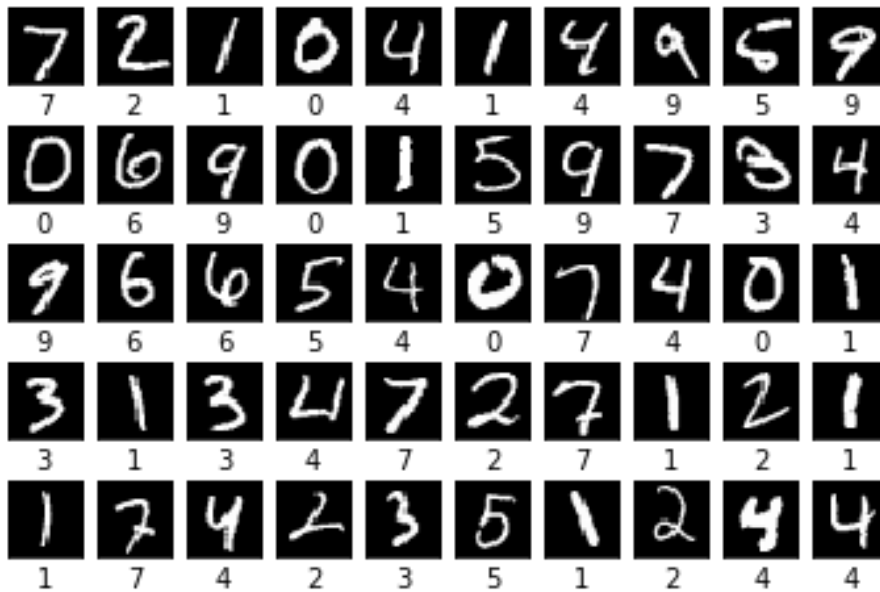
## Bouns 1

```
39. In [4]:
40. ### Bounus 1: show the first 40 images
41. ### ALways a good idea to validate that the data appears as you expect
42. ###
43. ### for sci-kit learn the images are represented as vectors of 784 elements
44. ### currently scaled from 0 to 1
45.
46. for i in range(50):
47.     l1_plot = plt.subplot(5, 10, i + 1)
```

```
48.     l1_plot.imshow(255 * X_test[i].reshape(28, 28), interpolation='nearest',
49.                  cmap=plt.cm.gray)
50.     l1_plot.set_xticks(())
51.     l1_plot.set_yticks(())
52.     #l1_plot.set_xlabel('Class %s' % y_test[i].decode())
53.     l1_plot.set_xlabel('%i' % int(y_test[i]))
54. plt.suptitle('Test image Examples')
55. plt.show()
```



Test image Examples

```
1.  ## Data standardization
2.  ## by the mean and standared deviation of the training set
3.  scaler = StandardScaler()
4.  X_train = scaler.fit_transform(X_train)
5.  X_test = scaler.transform(X_test)
```

```
6.  In [7]:
7.  #train and test classifier
8.  # Turn up tolerance for faster convergence
9.  clf = LogisticRegression(C=50. / 1000,
10.                  multi_class='multinomial',
11.                  penalty='l1', solver='saga', tol=0.1)
12. # Train the classifier
13. clf.fit(X_train, y_train)
14.
15. #Evaluate the classifier
16. sparsity = np.mean(clf.coef_ == 0) * 100
17. score = clf.score(X_test, y_test)
18. # print('Best C % .4f' % clf.C_)
19. print("Sparsity with L1 penalty: %.2f%%" % sparsity)
20. print("Test score with L1 penalty: %.4f" % score)
```

Sparsity with L1 penalty: 16.35%
Test score with L1 penalty: 0.9023

```python
1.  """ Custom datatset loader
2.      based on https://github.com/utkuozbulak/pytorch-custom-dataset-examples
3.  """
4.  import pandas as pd
5.  import imageio
6.
7.  class SimpleDataset():
8.      def __init__(self, data_path, csv_name, transform = None ):
9.          """
10.         Args:
11.             data_path (string): path to the folder where images and csv files are l
    ocated
12.             csv_name (string): name of the csv lablel file
13.             transform: pytorch transforms for transforms and tensor conversion
14.         """
15.         # Set path
16.         self.data_path = data_path
17.         # Read the csv file
18.         self.data_info = pd.read_csv(data_path + csv_name, header=None)
19.         # First column contains the image paths
20.         self.image_arr = np.asarray(self.data_info.iloc[:, 0])
21.         # Second column is the labels
22.         self.label_arr = np.asarray(self.data_info.iloc[:, 1])
23.         # Calculate len
24.         self.data_len = len(self.data_info.index)
25.
26.     def __getitem__(self, index):
27.         # Get image name from the pandas df
28.         single_image_name = self.image_arr[index]
29.         # Open image
30.         img_as_img = imageio.imread(self.data_path + single_image_name)
31.
32.         # Get label(class) of the image based on the cropped pandas column
33.         single_image_label = self.label_arr[index]
34.
35.         return (img_as_img, single_image_label)
36.
37.     def __len__(self):
38.         return self.data_len
```

```python
39. In [33]:
40. mydata = SimpleDataset( "./myDataSet/", "label.csv")
41.
42. #splitting into images and labels
43. X = []
44. y = []
45. for i in range(len(mydata)):
46.     X.append(mydata[i][0])
47.     y.append((mydata[i][1]))
48.
49. #converting into numpy arrays to enable easy reshaping and other array operations
50.
51. X = np.asarray(X)
52. print("Shape of the input image", X.shape)
53. y= np.asarray(y)
```

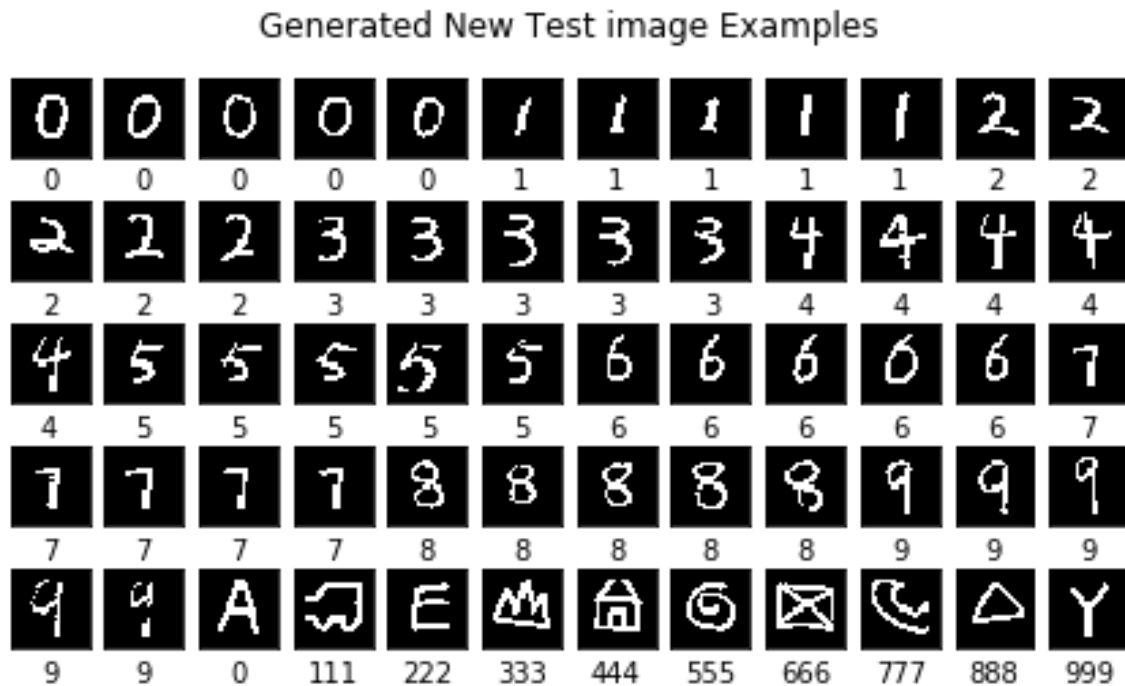Shape of the input image (60, 28, 28)

```python
1.  import warnings
```

```
2.  warnings.filterwarnings('ignore')
3.  #l1_plot = plt.subplot(12, 5, i + 1)
4.  for i in range(60):
5.      l1_plot = plt.subplot(5, 12, i + 1)
6.      l1_plot.imshow(X[i], interpolation='nearest',
7.                     cmap=plt.cm.gray)
8.      l1_plot.set_xticks(())
9.      l1_plot.set_yticks(())
10.     l1_plot.set_xlabel('%i' % int(y[i]))
11. plt.suptitle('Generated New Test image Examples')
12. plt.show()
```



Generated New Test image Examples

Note: above code are copied from my jupyter notebook.

# 5. Experiments with Logistic Regression classifier

Starting with the classifier described in the Lab tutorial evaluate the performance of the classifier for both the standard test dataset and also your own dataset.
Plot the images from the standard MNIST test set of the first 40 images with errors together with the correct and incorrect labels. Discuss these results.
Plot all your test images with truth and predicted labels. Discuss these results.

```
1.  #reshaping the array into flattened 784 array as an input for prediciton by the log
    istic regression classifier
2.  X = X.reshape(X.shape[0], 784)
3.  X = X / 255.
4.  #data standardiation with the training set statistics is required for this clasifie
    r
5.  X = scaler.transform(X)
6.
```

```
7.  y_pred = clf.predict(X)
8.
9.  score = clf.score(X, y)
10.
11. print("Test score with L1 penalty: %.4f" % score)
12. print("y_predicted_values", y_pred)
13. print("y_labels", y)
```

Test score with L1 penalty: 0.3000
y_predicted_values [9 9 9 9 9 1 1 1 1 1 8 8 9 2 2 3 3 3 3 3 3 8 9 8 8 8 8 8 9 8 8 8 8 8 6 8 1 1
 7 7 7 3 9 8 3 8 8 8 8 3 8 8 2 8 8 8 8 3 6 0 8]
y_labels [ 0  0  0  0  0  1  1  1  1  1  2  2  2  2  2  3  3  3
  3  3  4  4  4  4  4  5  5  5  5  5  6  6  6  6  6  7
  7  7  7  7  8  8  8  8  8  9  9  9  9  9  0 111 222 333
 444 555 666 777 888 999]

## Discussion:

The penalty = 0.3, it is quite high and good result, the quality of my data set is high, because I use Mark pen to draw every number and character, the stroke is thick and clear, even if we scale the size of image, the outline if every number is still clear, and we use scanner to scan our image into the computer, so the background and forground is very clear even before we do threthold operation to divide the forground and background.

I also write some interesting characters such as "telephone" "mail" "house"…. At first I want to label them as its name using string datastructure, but I find in the tutorial code sample, the data structure of the label is stricted to integer, so I label them as "000" "111"…"999", in python, "000" = "0", so from the result, the label "000" is converted to "0" automatically.

Then I put out the confusion matrix.

```
1.  ## For analysis show also the confusion matrix
2.
3.  from sklearn.metrics import confusion_matrix
4.  y_predict = clf.predict(X_test)
5.  cfm = confusion_matrix(y_test, y_predict)
6.  print (cfm)
```

```
[[ 957    0    1    1    0    8    9    2    2    0]
 [   0 1104    2    3    1    2    4    0   19    0]
 [  15   16  888   24   16    1   12   21   35    4]
 [   5    4   19  901    1   32    6   17   17    8]
 [   2    8    4    0  916    2   12    2    3   33]
 [  15    5    4   41   17  743   20   14   25    8]
 [  16    6    7    0   12   19  894    1    3    0]
 [   1   25   18    5   12    0    1  928    1   37]
 [  11   23   10   26   19   33   11   17  810   14]
 [  16    8    3   13   48    7    0   28    4  882]]
```
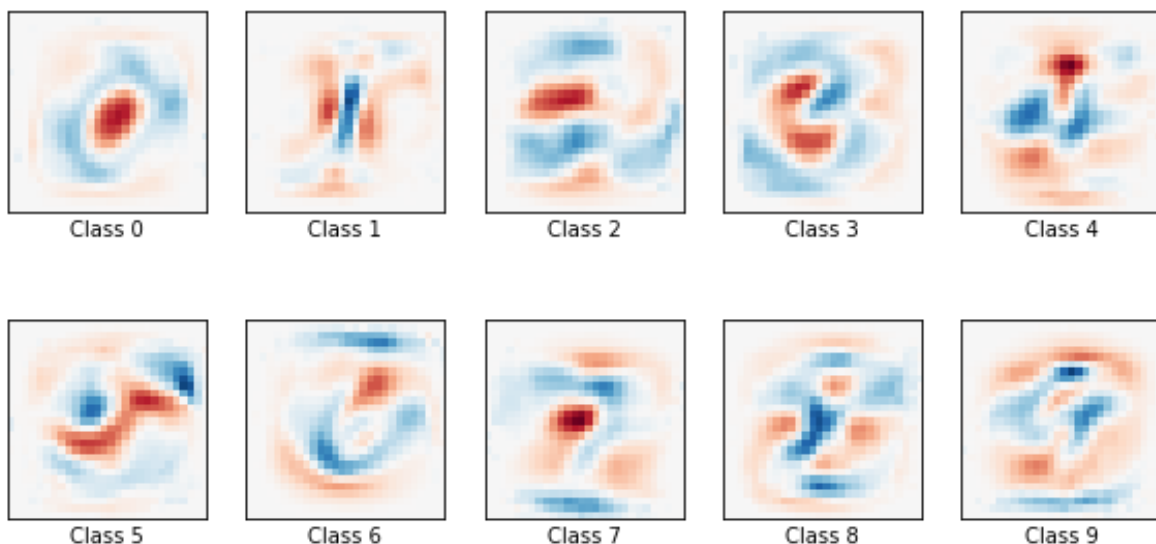
## Bonus 2

```
1.  # Bonus 2: Visualization of the weights
2.  # This is only possible for simple classifiers
3.
4.  coef = clf.coef_.copy()
5.  plt.figure(figsize=(10, 5))
6.  scale = np.abs(coef).max()
7.  for i in range(10):
8.      l1_plot = plt.subplot(2, 5, i + 1)
9.      l1_plot.imshow(coef[i].reshape(28, 28), interpolation='nearest',
10.                 cmap=plt.cm.RdBu, vmin=-scale, vmax=scale)
11.      l1_plot.set_xticks(())
12.      l1_plot.set_yticks(())
13.      l1_plot.set_xlabel('Class %i' % i)
14.  plt.suptitle('Classification vector for...')
15.  plt.show()
```

Classification vector for...



Class 0    Class 1    Class 2    Class 3    Class 4

Class 5    Class 6    Class 7    Class 8    Class 9

## Discussion:

The weight of visualization is quite straight forward, the orange color shows the contour of every forground, look at the orange colorin every class, class 0 looks like the number 0, class 1 looks like 1…. But for class 5, the orange color is more obscure, it is hard for us to imagine the contour of "5" from the orange color. Maybe "5" is more complicated in shape compared with other numbers.

## Bouns 3

```
1.  ### Bonus 3: comfusion matrix visualization tool
2.  ## A more elegant preserntation for a confusion matrix
3.
```

```
4.   import itertools
5.   def plot_confusion_matrix(cm, classes,
6.                             normalize=False,
7.                             title='Confusion matrix',
8.                             cmap=plt.cm.Blues):
9.       """
10.      This function prints and plots the confusion matrix.
11.      Normalization can be applied by setting `normalize=True`.
12.      """
13.      if normalize:
14.          cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
15.          print("Normalized confusion matrix")
16.      else:
17.          print('Confusion matrix, without normalization')
18.
19.      #print(cm)
20.
21.      plt.imshow(cm, interpolation='nearest', cmap=cmap)
22.      plt.title(title)
23.      plt.colorbar()
24.      tick_marks = np.arange(len(classes))
25.      plt.xticks(tick_marks, classes, rotation=45)
26.      plt.yticks(tick_marks, classes)
27.
28.      fmt = '.2f' if normalize else 'd'
29.      thresh = cm.max() / 2.
30.      for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
31.          plt.text(j, i, format(cm[i, j], fmt),
32.                   horizontalalignment="center",
33.                   color="white" if cm[i, j] > thresh else "black")
34.
35.      plt.ylabel('True label')
36.      plt.xlabel('Predicted label')
37.      #plt.tight_layout()
38.      plt.show()
39.  plot_confusion_matrix(cfm, classes=range(10), normalize=True,
40.                        title='Confusion matrix for MNIST')
```
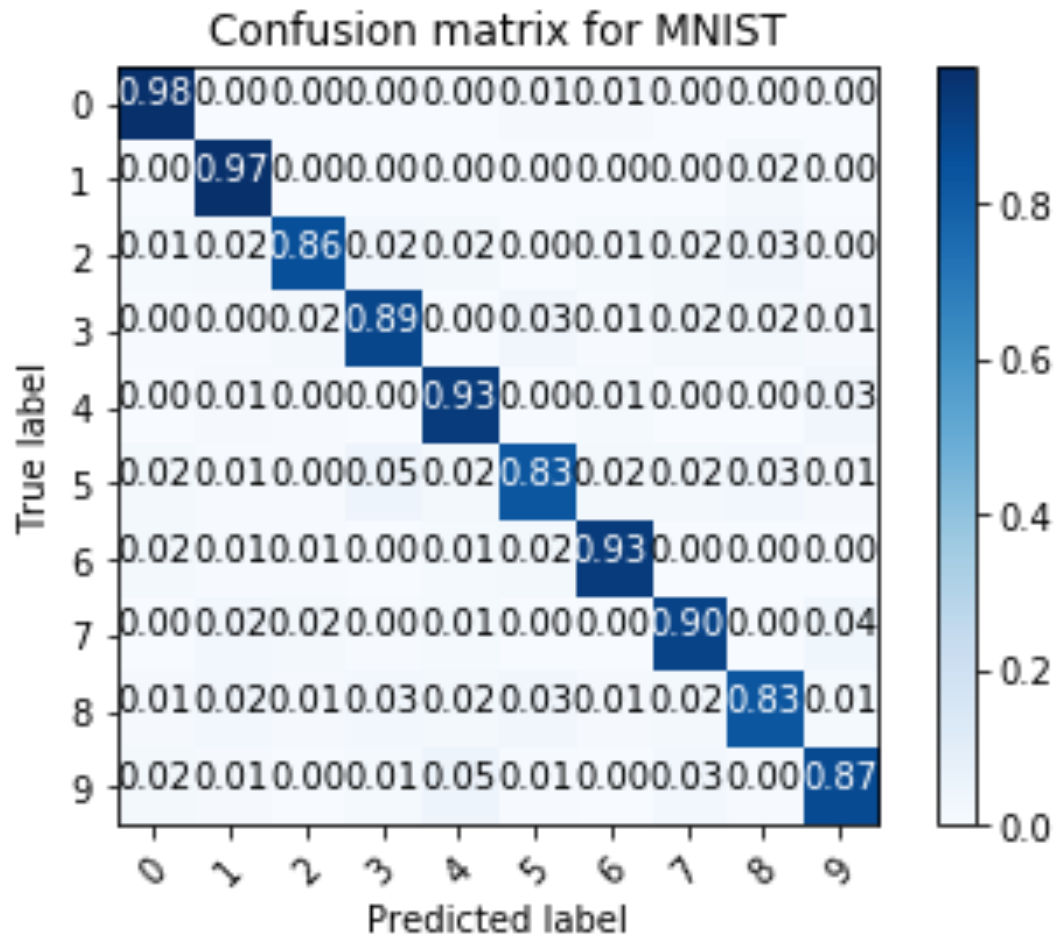
Finally, we plot the confusion matrix for MNIST as below:

Normalized confusion matrix

Confusion matrix for MNIST

## Discussion：

The sample code in tutorial is not absolutely correct, I use the code from tutorial cannot draw the confusion matrix for MNIST, I make small change on the code, as below:
I add

```
1.  normalize=True,
```

I comment this line and add another plot line.

```
1.  #plt.tight_layout()
2.      plt.show()
```