

Machine Learning HW3 – Neural Network

NCTU EE 0310128 游騰杰

Programming language: Python

Environment: Python 2.7.13 under Homebrew @ Mac OS X 10.11.6

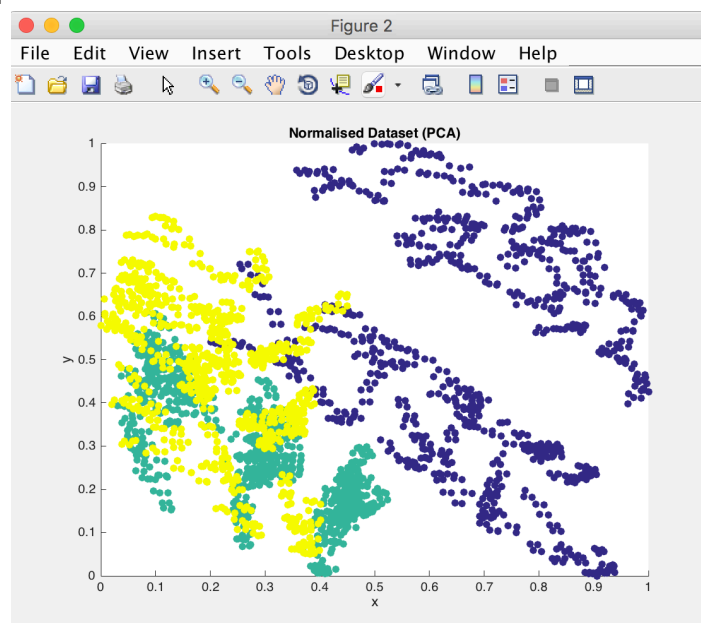
Python module used: numpy, pandas, math, Pillow, random, scikit-learn

Plot: Matlab

1. Use Principal component analysis (PCA) to map data down to 2 dimensions.

```
21
22 def dimension_reduction(data, test, n, bLDA, target):
23     if not bLDA:
24         dimReduction = sklearn.PCA(n_components=n)
25         transformed_data = dimReduction.fit_transform(data)
26         transformed_test = dimReduction.transform(test)
27     else:
28         dimReduction = sklearn.LDA(n_components=n)
29         transformed_data = dimReduction.fit_transform(data, target)
30         transformed_test = dimReduction.transform(test)
31
32     return transformed_data, transformed_test
33
51 # Dimension Reduction
52 transformed_data, transformed_test \
53 = nn.dimension_reduction(cv_data, testing_data, reduce_dimension, \
54                          bLDA, train_class)
55 #nn.exportData3d("pca.csv", transformed_data[:, 0] \
56 #                , transformed_data[:, 1] \
57 #                , train_class)
```

Dataset in feature subspace:

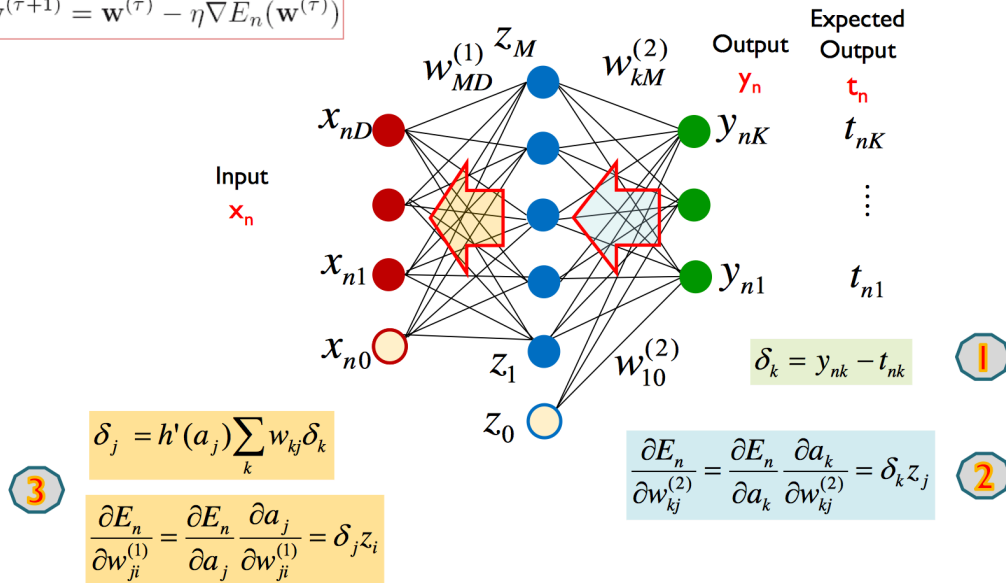


Reference: How to Implement the Backpropagation Algorithm From Scratch In Python

<http://machinelearningmastery.com/implement-backpropagation-algorithm-scratch-python/>

2. Use stochastic gradient descent in back propagation.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})$$



Reference: (Lecture Notes) Neural Networks P. 21

I begin calculating all the deltas in every neuron backward-wise, starting with the output layer prediction with formula

$$\delta_k = y_k - t_k \quad (5.54)$$

[Formula] Reference: Textbook P. 243 (5.54)

all the preceding layers have the delta of the form

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k \quad (5.56)$$

[Formula] Reference: Textbook P. 244 (5.56)

and I can obtain all the deltas, next combining those delta with the learning rate (η) and the layer inputs z_i' s or z_j' s using the Stochastic Gradient Descent formula to update the weight in each iteration.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)}). \quad (5.43)$$

[Formula] Reference: Textbook P. 240 (5.43)

Implementation:

```

305
306 def training_nn(network, data, target, l_rate, n_epoch, bReLU, bMiniBatch, batch_iter):
307     # target shape: numData x numClass (in one-of-k format)
308     numData, numClass = target.shape
309     numLayer = len(network)
310     batch_count = int(numData / batch_iter)
311     iterations = batch_iter if bMiniBatch else numData
312     # Randomise the sequence of data
313     bundle = np.random.permutation(np.hstack([data, target]))
314     suffle_data = bundle[:, :data.shape[1]]
315     suffle_targ = bundle[:, -1 * numClass:]
316     #
317     for epoch in range(n_epoch):
318         error = np.array([0.0] * n_epoch)
319         if not bMiniBatch:
320             for itr_Data in range(numData):
321                 reset_delta_gradE(network)
322                 #probability = forward_propagate(network, data[itr_Data], bReLU)
323                 probability = forward_propagate(network, suffle_data[itr_Data], bReLU)
324                 prediction = softmax(probability)
325                 # Current class target
326                 #curr_target = target[itr_Data, :]
327                 curr_target = suffle_targ[itr_Data, :]
328                 error[epoch] = sum([pow(curr_target[itr] - prediction[itr], 2) \
329                                     for itr in range(numClass)])
330                 calError_backprop(network, curr_target, bReLU)
331                 update_weights(network, suffle_data[itr_Data], l_rate, bMiniBatch, False)
332

```

Result:

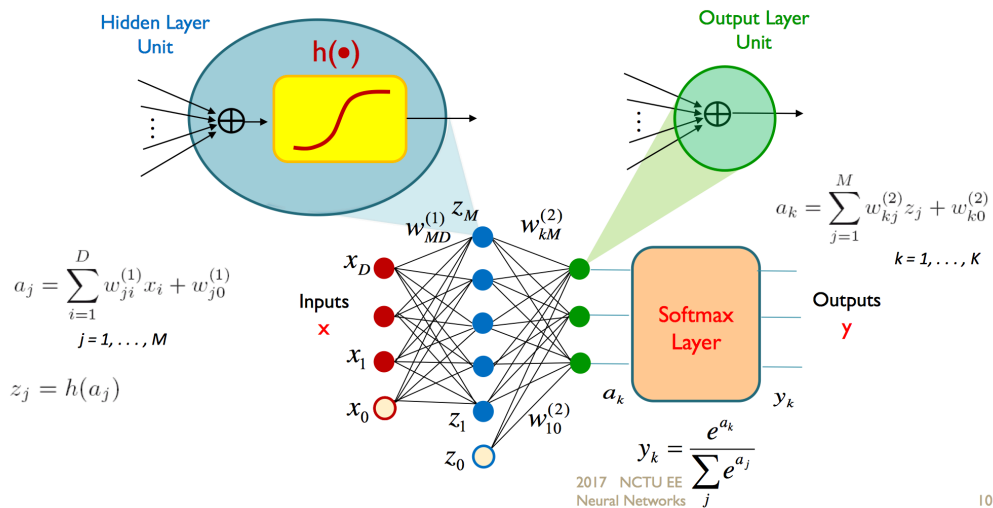
1 Hidden Layer	2 Hidden Layers
<p>epoch: 499, lrate: [0.1, 0.1], error: 0.383 epoch: 500, lrate: [0.1, 0.1], error: 0.383</p> <p>==== Summary =====</p> <p>Neural Network:</p> <ul style="list-style-type: none"> - Input Layer : 2 neurons - Hidden Layer #1: 4 neurons - Output Layer : 3 neurons <p>Error rate: 11.50% (69 of 600)</p> <p>In-Class error rate:</p> <ul style="list-style-type: none"> - Class 1: 9.00% (18 of 200) - Class 2: 0.50% (1 of 200) - Class 3: 25.00% (50 of 200) 	<p>epoch: 200, lrate: [0.1, 0.1, 0.1], error: 0.384</p> <p>==== Summary =====</p> <p>Neural Network:</p> <ul style="list-style-type: none"> - Input Layer : 2 neurons - Hidden Layer #1: 5 neurons - Hidden Layer #2: 7 neurons - Output Layer : 3 neurons <p>Error rate: 12.00% (72 of 600)</p> <p>In-Class error rate:</p> <ul style="list-style-type: none"> - Class 1: 9.50% (19 of 200) - Class 2: 3.50% (7 of 200) - Class 3: 23.00% (46 of 200)
1 Hidden Layer (Cross-Validation)	2 Hidden Layers (Cross-Validation)
<p>epoch: 500, lrate: [0.1, 0.1], error: 0.419 Iteration: #4 Error rate: 16.83% (101 of 600)</p> <p>=== Cross-Validation Summary ===</p> <p>Iterations : 4 Epoch : 500 Average error rate: 14.58% Maximum error rate: 16.83% Minimum error rate: 10.83%</p>	<p>epoch: 199, lrate: [0.1, 0.1, 0.1], error: 0.434 epoch: 200, lrate: [0.1, 0.1, 0.1], error: 0.437 Iteration: #4 Error rate: 25.83% (155 of 600)</p> <p>=== Cross-Validation Summary ===</p> <p>Iterations : 4 Epoch : 200 Average error rate: 16.54% Maximum error rate: 25.83% Minimum error rate: 11.83%</p>

3. Implement the neural network model with 1-hidden layer. Choose the sigmoid function as the activation function.

Design Structure:

Feed-forward Neural Network (9/10)

- Multi-Class Classification Model



Reference: (Lecture Notes) Neural Networks P. 10

First connect the inputs (red nodes) to the hidden layer (blue nodes) using the formula

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \quad (5.2)$$

[Formula] Reference: Textbook P. 227 (5.2)

where a_j denotes the “activation” of the node, in addition, I choose sigmoid function as the activation function (shown below, equation 5.5)

```

168
169 def sigmoid_transferFcn(activation):
170     # Logistic Sigmoid function
171     return 1.0 / (1.0 + math.exp(-activation))
172
173 def sigmoid_derivative(x):
174     # sigmoid function derivative = output * (1.0 - output)
175     return x * (1.0 - x)
176

```

$$y_k = \sigma(a_k) \quad (5.5)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}. \quad (5.6)$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma). \quad (4.88)$$

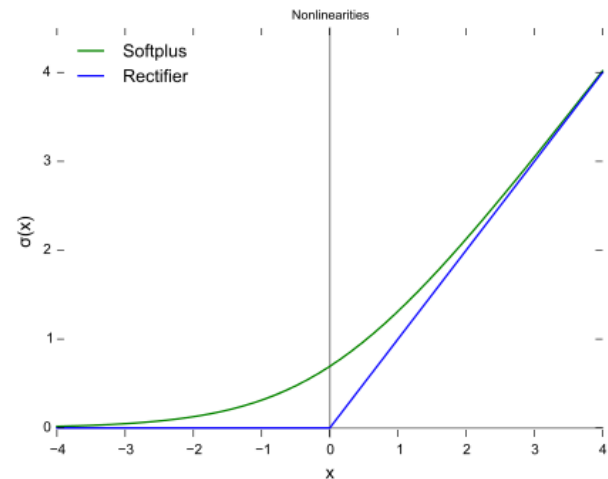
[Formula] Reference: Textbook P. 205 (4.88), P. 228 (5.5) and (5.6)

Continue this process in the output layer (green nodes), but here the sigmoid function is changed to Softmax function for the multi-class classification model.

4. Implement the neural network model with 2-hidden layer. Choose the rectified function as the activation function.

```
176
177 def rectify_transferFcn(activation):
178     # Rectify function
179     alpha = -0.01
180     return activation if activation > 0.0 else \
181         alpha * activation
182
183 def rectify_derivative(x):
184     # Rectify function = 0 if x <= 0 else 1
185     alpha = -0.01
186     #print x
187     return 1. if x > 0.0 else alpha
188
```

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$



Reference: Rectifier (neural networks) - Wikipedia

[https://en.wikipedia.org/wiki/Rectifier_\(neural_networks\)](https://en.wikipedia.org/wiki/Rectifier_(neural_networks))

The process and structure is basically the same as the previous one, but the number of hidden layer is added to two, also the activation function, which is sigmoid function in the previous case, is implemented with rectify function, or rectify linear unit (ReLU) function.

At first I found it difficult to train the model with the rectify function, after searching for some possible causes, it turned out that the rectify region ($x < 0$) might suffer from the dying ReLU problem in some applications, so here I use leaky ReLU for rectify function instead.

Leaky ReLU. Leaky ReLUs are one attempt to fix the “dying ReLU” problem. Instead of the function being zero when $x < 0$, a leaky ReLU will instead have a small negative slope (of 0.01, or so). That is, the function computes $f(x) = \mathbb{1}(x < 0)(\alpha x) + \mathbb{1}(x \geq 0)(x)$ where α is a small constant. Some people report success with this form of activation function, but the results are not always consistent. The slope in the negative region can also be made into a parameter of each neuron, as seen in PReLU neurons, introduced in [Delving Deep into Rectifiers](#), by Kaiming He et al., 2015. However, the consistency of the benefit across tasks is presently unclear.

Reference: CS231n Convolutional Neural Networks for Visual Recognition

<http://cs231n.github.io/neural-networks-1/>

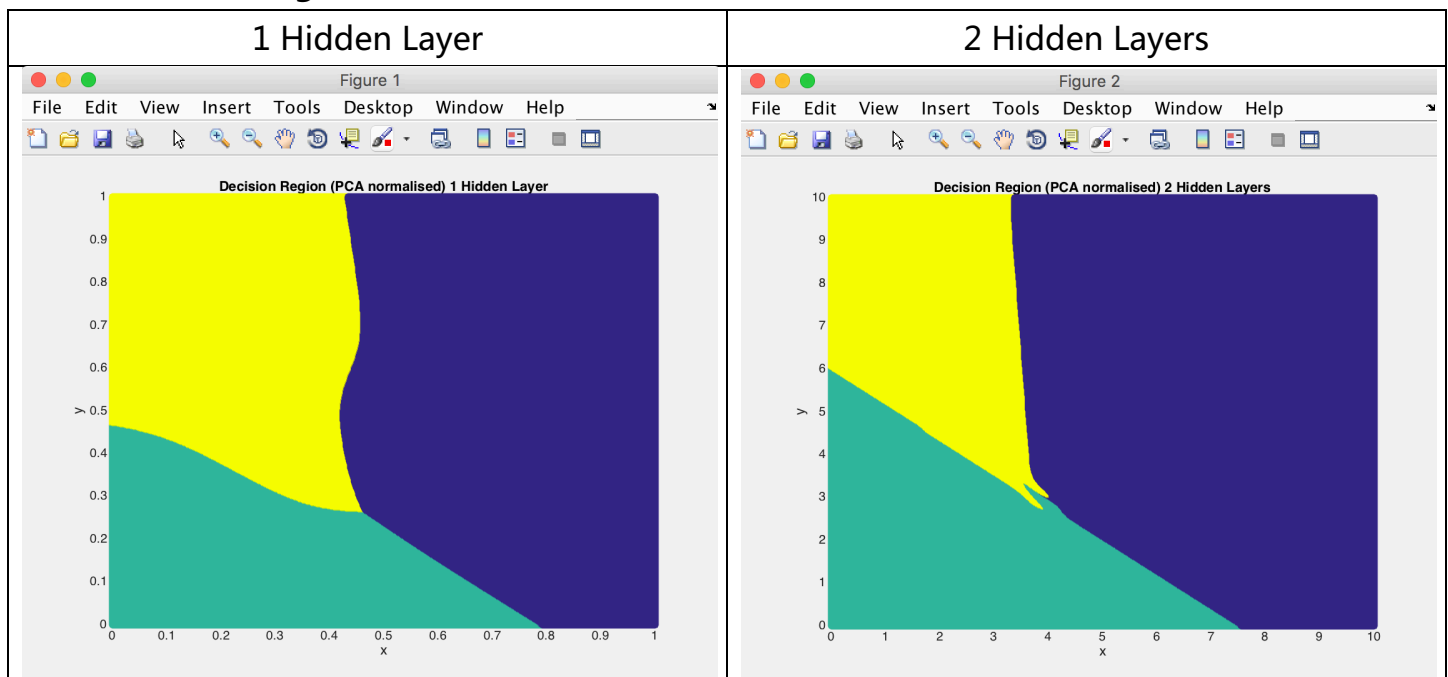
Implementing on Python:

```
7 # Constants / Parameter
8 cv_part = 3
9 cv_enable = False
10 cv_count = cv_part if cv_enable else 1
11
12 suffle = True
13 reduce_dimension = 2
14
15 bLDA = False
16
17 hidden = 5
18 bReLU = False
19
20 bMiniBatch = True
21 batch_size = 5
22 batch_iter = batch_size if bMiniBatch else 1
23
24 reserved = [0.20, 0.20, 0.20]
25 l_rate = 0.10
26 n_epoch = 200
27 scaling = 1.0
28
29 Train_Path = "Data_Train/"
30 Demo_Path = "Demo/"
31
```

```
7 # Constants / Parameter
8 cv_part = 20
9 cv_enable = False
10 cv_count = cv_part if cv_enable else 1
11
12 suffle = True
13 reduce_dimension = 2
14
15 bLDA = False
16
17 hidden_1 = 5
18 hidden_2 = 7
19 bReLU = False
20
21 bMiniBatch = False
22 batch_size = 50
23 batch_iter = batch_size if bMiniBatch else 1
24
25 reserved = [0.10, 0.10, 0.10]
26 l_rate = 0.05
27 n_epoch = 25
28 scaling = 10.0
29
30 Train_Path = "Data_Train/"
31 Demo_Path = "Demo/"
32
```

Name	Function
cv_part	Determine the amount of n in n-fold cross validation.
cv_enable	Enable or disable the cross-validation functioning in the code.
cv_count	Determine the times of iteration based on previous two values.
suffle	Determine whether to shuffle the input data sequence.
reduce_dimension	Determine the dimension of PCA reduces to.
bLDA	Determine whether use LDA instead of PCA.
hidden/hidden1/ hidden2	Set the number of neurons in the hidden layers.
bReLU	Determine whether to use Rectify Linear Function (ReLU) as activation function.
bMiniBatch	Determine whether to use mini-Batch Gradient Descent.
batch_size	Set the size of each mini-Batch.
batch_iter	Set the size of each mini-Batch according to bMiniBatch.
reserved_part	Define the portion of the data for each class that saved for testing (the number * 100 equal percentage of given data).
l_rate	Learning rate.
n_epoch	Set the number of epochs.
Train_Path	A list of iterations that a w of each class required to converge.
Demo_Path	The order of polynomial basis function.

5. Plot decision regions.



6. Compare and discuss their performances with homework 2.

Comparison:

	Linear Model		Neural Networks		
	Generative	Discriminative	SGD w/ sigmoid (1 hidden layer)	SGD w/ ReLU (2 hidden layers)	mBGD w/ sigmoid (1 hidden layer)
Error Rate	18.33%	11.00%	11.50%	12.00%	12.33%

SGD: Stochastic Gradient Descent

mBGD: mini-Batch Gradient Descent

ReLU: Rectify Linear Unit function

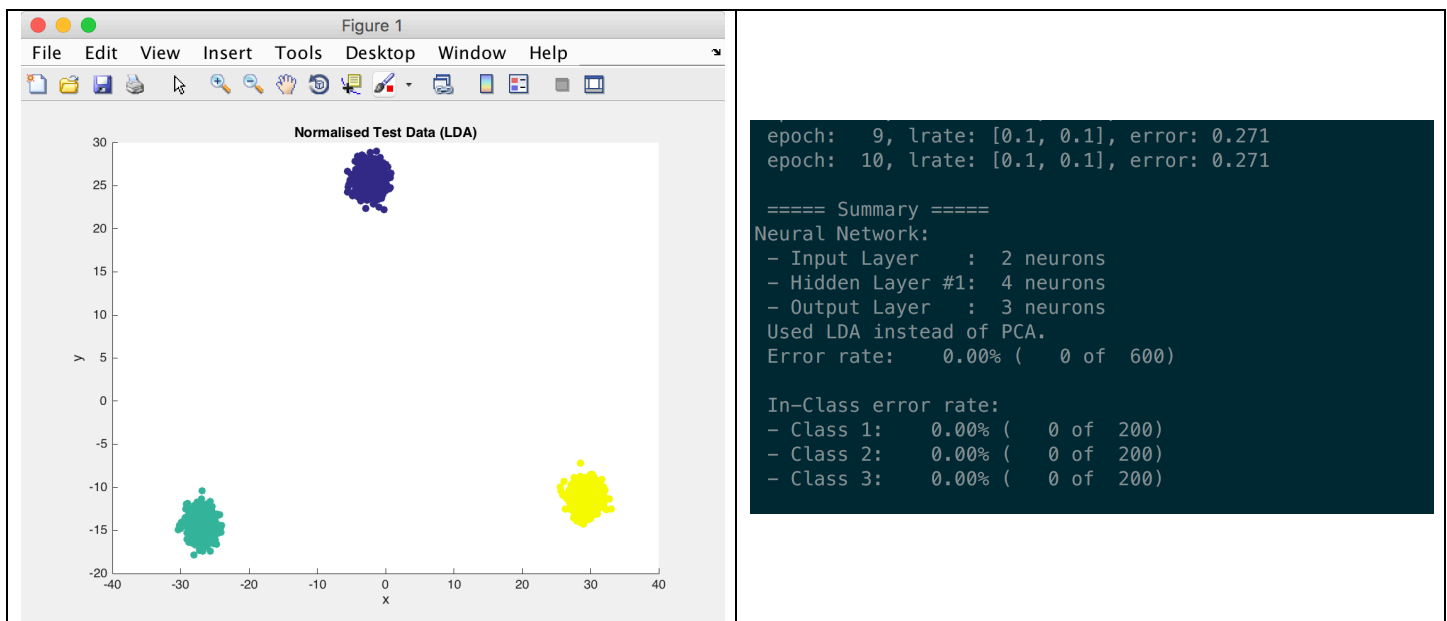
Discussion:

Starting from the training time between these models, both of the Linear Models have the least training time and the result is moderate; Neural Networks are models with higher computational complexity, but the performance and Cross-Validation result showed of high stability in the performance and high accuracy. However, the difference is not way too big, perhaps the problem itself is not complicated enough for the Neural Networks to stand out from those models, also, increasing the number of neurons in each layer may result in better fitting to the data.

Mini-Batch Gradient Descent (mBGD) is aimed to eliminate the variance between consecutive data, which leads to more stable change in the process of Gradient Descent, in the case of Stochastic Gradient Descent (SGD), the model would be too sensitive to the outliers.

Determining the size of each mini-Batch is of minor importance, I tried several times and found that 5 is a moderate number, since my training data is of a random class sequence, so if size of the mini-Batch is too large, the variance of the mini-Batch would be too small for training, or perhaps more epochs is needed. In my case, a mini-Batch with size of 5 entries of data converges quickly with around 50 epochs.

Additionally, I found the PCA is unable to seek the maximum separation between classes, so I tried performing Linear Discriminant Analysis (LDA) on the dataset, which is frequently used in supervised learning, the distribution of the dataset is shown below, both the distribution of the dataset and the result is outstanding.



7. Explain and compare the following nouns with words:

(1) Batch Gradient Descent (BGD):

Use entire the dataset as training set, take average of the gradient of error function within the entire training set, and then update the weights once each epoch.

$$\mathbf{w}^{(old)} = \mathbf{w}^{new} + \eta \cdot \nabla E_{avg \text{ (whole trainingset)}}(\mathbf{w})$$

(2) mini-Batch Gradient Descent (mBGD):

Divide the entire dataset to smaller dataset, take average of the gradient of error function within the mini-batch dataset, and then update the weights according to the average.

$$\mathbf{w}^{(old)} = \mathbf{w}^{new} + \eta \cdot \nabla E_{avg \text{ (mini-batch)}}(\mathbf{w})$$

(3) Stochastic Gradient Descent (SGD):

Calculate the gradient of error function for each data in the training set and update the weights after each propagation of data.

$$\mathbf{w}^{(old)} = \mathbf{w}^{new} + \eta \cdot \nabla E_{stochastic}(\mathbf{w})$$

(4) Online Gradient Descent (OGD):

Calculate the gradient of error function after gaining real-time data and update the weights after the data propagates through the network. The big difference is that OGD discards the all the data after training.

$$\mathbf{w}^{(old)} = \mathbf{w}^{new} + \eta \cdot \nabla E_{online}(\mathbf{w})$$

OGD is used mostly in applications involving the analysis the trend of a phenomenon, which requires real-time accuracy, collecting data to a certain amount and update the model would be impractical for some applications, e.g. analyzing the click through rate (CTR) of an advertisement, if the data is collected (batch learning) and update the model once a fixed interval, say once a day, it would be relatively time-costing and the model cannot reflect the real-time trend of the CTR.

Comparison:

	BGD	mBGD	SGD	OGD
Training set	Fixed	Fixed	Fixed	Real-time update
Data in each iteration	Entire dataset	Subsets from the entire dataset	One entry in the entire dataset	Depends
Required computational power	High	Middle	Low	Low
Convergence (Stability)	Most stable	Comparatively stable	Not stable	Not stable

8. Bonus:

Use mini-batch method to retry the task 3 and add discussion into task 6.

```
333         else:
334             offset = 0
335             error = [0.0] * n_epoch
336             mbatch_data, mbatch_targ \
337                 = mini_batch_training(suffle_data, suffle_targ, batch_iter)
338             for itr_batch in range(batch_count):
339                 reset_delta_gradE(network)
340                 # Load current batch data and target
341                 b_data = mbatch_data[itr_batch]
342                 b_targ = mbatch_targ[itr_batch]
343                 for itr_Data in range(batch_iter):
344                     probability = forward_propagate(network, b_data[itr_Data], bReLU)
345                     prediction = softmax(probability)
346                     # Current class target
347                     curr_target = b_targ[itr_Data]
348                     error[epoch] += sum([pow(curr_target[itr] - prediction[itr], 2) \
349                                         for itr in range(numClass)])
350                     calError_backprop(network, curr_target, bReLU)
351                     # miniBatch update grad(E)
352                     update_weights(network, b_data[itr_Data], l_rate, bMiniBatch, True)
353
354                 # miniBatch update weights
355                 update_weights(network, b_data[0, :], l_rate, bMiniBatch, False)
356
357             error[epoch] /= (batch_iter * batch_count)
358
359             print " epoch: {:3d}, lrate: {:.3f}, error: {:.3f}" \
360                   .format(epoch + 1, l_rate, error[epoch])
361
246 def update_weights(network, data, l_rate, bMiniBatch, updateGradE):
247     # w(new) = w(old) + learning_rate * delta * input
248     # mini-batch: w(new) = w(old) + learning_rate * avg(grad(E(w)))
249     numLayer = len(network)
250     for itr_layer in range(numLayer):
251         if itr_layer != 0:
252             # Inputs of current layer = outputs from previous layer
253             inputs = [neuron['output'] for neuron in network[itr_layer - 1]]
254         else:
255             # Initialize input with input data (bias not included)
256             inputs = data
257
258         numInput = len(inputs)
259         for neuron in network[itr_layer]:
260             for itr in range(numInput): # bias is added later
261                 if not bMiniBatch:
262                     neuron['weights'][itr] += l_rate * neuron['delta'] * inputs[itr]
263                 elif updateGradE:
264                     neuron['gradE'][itr] = np.append(neuron['gradE'][itr],
265                                                       neuron['delta'] * inputs[itr])
266                 else:
267                     neuron['weights'][itr] += l_rate * neuron['gradE'][itr].mean()
268
269         # Update bias weight / gradient E
270         if not bMiniBatch:
271             neuron['weights'][-1] += l_rate * neuron['delta'] * 1.0
272         elif updateGradE:
273             neuron['gradE'][-1] = np.append(neuron['gradE'][-1],
274                                              neuron['delta'] * 1.0)
275         else:
276             neuron['weights'][-1] += l_rate * neuron['gradE'][-1].mean()
277
```

Result:

```
epoch: 195, lrate: 0.100, error: 0.433
epoch: 196, lrate: 0.100, error: 0.433
epoch: 197, lrate: 0.100, error: 0.433
epoch: 198, lrate: 0.100, error: 0.433
epoch: 199, lrate: 0.100, error: 0.433
epoch: 200, lrate: 0.100, error: 0.433

===== Summary =====
Activation function      : Sigmoid Function
mini-Batch Gradient Descent: True
mini-Batch Size         : 5
Error rate: 12.33% ( 74 of 600)

In-Class error rate:
- Class 1: 12.00% ( 24 of 200)
- Class 2:  2.00% (  4 of 200)
- Class 3: 23.00% ( 46 of 200)
```