

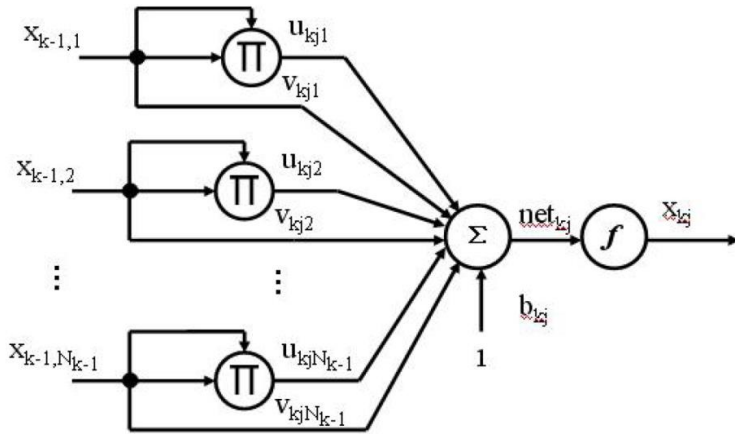
HW2 Report

姓名： 繆佳宇

学号: 51502191092

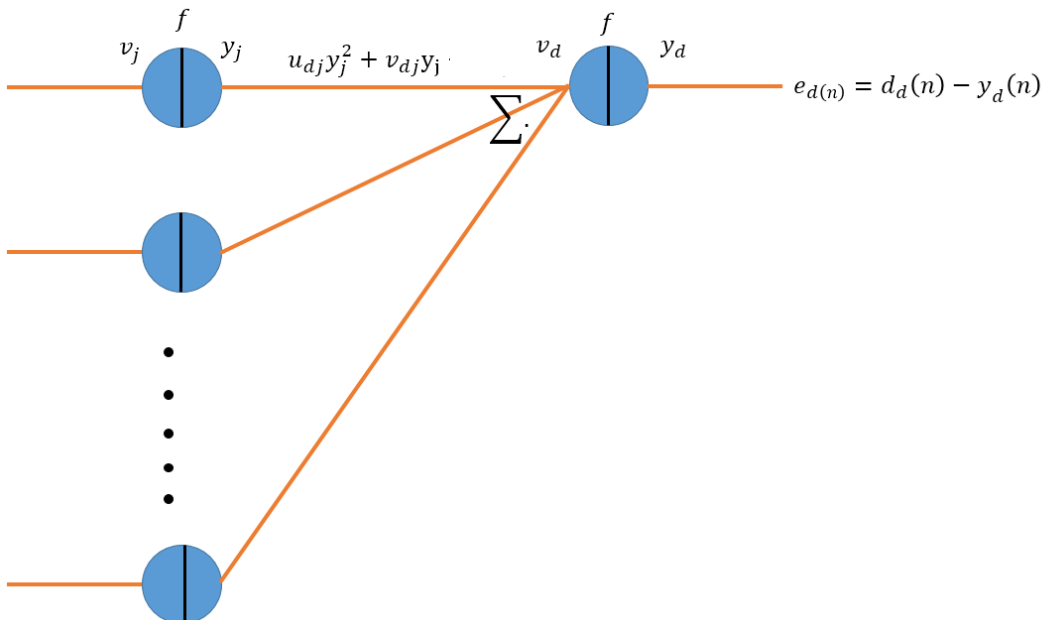
1. 试推导训练上述 MLQP 的下列两种 BP 算法:

a) 批量学习; b) 在线学习。



Solution:

先拿**在线学习**来说，MLQP 的 BP 算法和平常的 BP 的算法区别在于前向传播中神经元的输出 f 的作用域不是线性的而是二次的，但是推导是相近的。



输出层损失是 $\varepsilon(n) = \frac{1}{2} \sum_{k \in C} e_k^2(n)$ 。

如果节点 j 是输出层的节点，那么损失函数在节点 j 上的梯度

$$\begin{aligned}\delta_j(n) &= \frac{\partial \varepsilon(n)}{\partial v_j(n)} = \frac{\partial \varepsilon(n)}{\partial e_j(n)} \cdot \frac{\partial e_j(n)}{\partial y_j(n)} \cdot \frac{\partial y_j(n)}{\partial v_j(n)} \\ \frac{\partial \varepsilon(n)}{\partial e_j(n)} &= e_{j(n)} = d_j(n) - y_j(n) \\ \frac{\partial e_j(n)}{\partial y_j(n)} &= -1 \\ \frac{\partial y_j(n)}{\partial v_j(n)} &= \varphi'_j(v_j(n)) (\varphi \text{ 是激活函数}) \\ \delta_j(n) &= -e_{j(n)} \varphi'_j(v_j(n))\end{aligned}$$

如果节点 j 是中间节点，与之相连的下一层节点集合为 \mathbf{C} ，那么损失函数在节点 j 上的梯度

$$\begin{aligned}\delta_j(n) &= \frac{\partial \varepsilon(n)}{\partial v_j(n)} = \sum_{d \in \mathbf{C}} \frac{\partial \varepsilon(n)}{\partial v_d(n)} \cdot \frac{\partial v_d(n)}{\partial y_j(n)} \cdot \frac{\partial y_j(n)}{\partial v_j(n)} \\ \frac{\partial \varepsilon(n)}{\partial v_d(n)} &\text{ 即是 } \delta_d(n) \\ \frac{\partial v_d(n)}{\partial y_j(n)} &= 2u_{j,d}y_{j(n)} + v_{j,d} \\ \frac{\partial y_j(n)}{\partial v_j(n)} &= \varphi'_j(v_j(n)) (\varphi \text{ 是激活函数}) \\ \delta_j(n) &= \varphi'_j(v_j(n)) \sum_{d \in \mathbf{C}} \delta_d(n) (2u_{j,d}y_{j(n)} + v_{j,d})\end{aligned}$$

所以从最后一层节点反向传播误差，将各个节点的梯度都可以求出。

所以计算

$$\begin{aligned}\Delta u_{d,j} &= -\alpha \delta_d(n) y_j(n)^2, & \Delta v_{d,j} &= -\alpha \delta_d(n) y_j(n). \\ \Delta b_d &= -\alpha \delta_d(n) & (\alpha \text{ 是学习率})\end{aligned}$$

批量学习同理亦然，不同的是损失函数是批量样本的损失和。

$$\varepsilon(n) = \frac{1}{2N} \sum_n \sum_{k \in \mathbf{C}} e_k^2(n)$$

最后的参数更改也要相应改变。

$$\begin{aligned}\Delta u_{d,j} &= -\frac{\alpha}{N} \delta_d(n) y_j(n)^2, & \Delta v_{d,j} &= -\frac{\alpha}{N} \delta_d(n) y_j(n). (\alpha \text{ 是学习率}) \\ \Delta b_d &= -\frac{\alpha}{N} \delta_d(n)\end{aligned}$$

2.

编程实现问题 1) 中的训练 MLQP 的**在线 BP** 算法（编程语言不限，可采用常用的程序设计语言或 Matlab）

Solution:

这里我使用的是 python 来实现的，为了和第三题同步，首先我实现了一个带有一层隐层的神经网络的类，然后将 BP 算法作为它的一个方法，主要就是先得到后两层的节点梯度，然

后得到各个参数值的梯度下降值。（同时注意是在线 BP 算法）

程序截图如下：

```
101 def BackPropagate(self, x, y, lr):
102     """
103     the implementation of BP algorithms
104
105     @param x, y: array, input and output of the data sample
106     @param lr: float, the learning rate of gradient decent iteration
107     """
108
109     # dependent packages
110     import numpy as np
111
112     # get current network output
113     self.Pred(x)
114
115     # calculate the gradient based on output
116     o_grid = np.zeros(self.o_n)
117     for j in range(self.o_n):
118         o_grid[j] = (y[j] - self.o_v[j]) * self.afd(self.o_v[j])
119
120     h_grid = np.zeros(self.h_n)
121     for h in range(self.h_n):
122         for j in range(self.o_n):
123             h_grid[h] += (2*self.ho_u[h][j]*self.h_v[h]+self.ho_v[h][j]) * o_grid[j]
124         h_grid[h] = h_grid[h] * self.afd(self.h_v[h])
125
126     # updating the parameter
127     for h in range(self.h_n):
128         for j in range(self.o_n):
129             self.ho_u[h][j] += lr * o_grid[j] * (self.h_v[h]**2)
130             self.ho_v[h][j] += lr * o_grid[j] * self.h_v[h]
131
132     for i in range(self.i_n):
133         for h in range(self.h_n):
134             self.ih_u[i][h] += lr * h_grid[h] * (self.i_v[i]**2)
135             self.ih_v[i][h] += lr * h_grid[h] * self.i_v[i]
136
137     for j in range(self.o_n):
138         self.o_t[j] -= lr * o_grid[j]
139
140     for h in range(self.h_n):
141         self.h_t[h] -= lr * h_grid[h]
```

3.

用含有一层中间层（中间层的神经元数目可设成 10）的两层 MLQP 学习双螺旋问题，比较在三种不同学习率下网络的训练时间和决策面。

Solution:

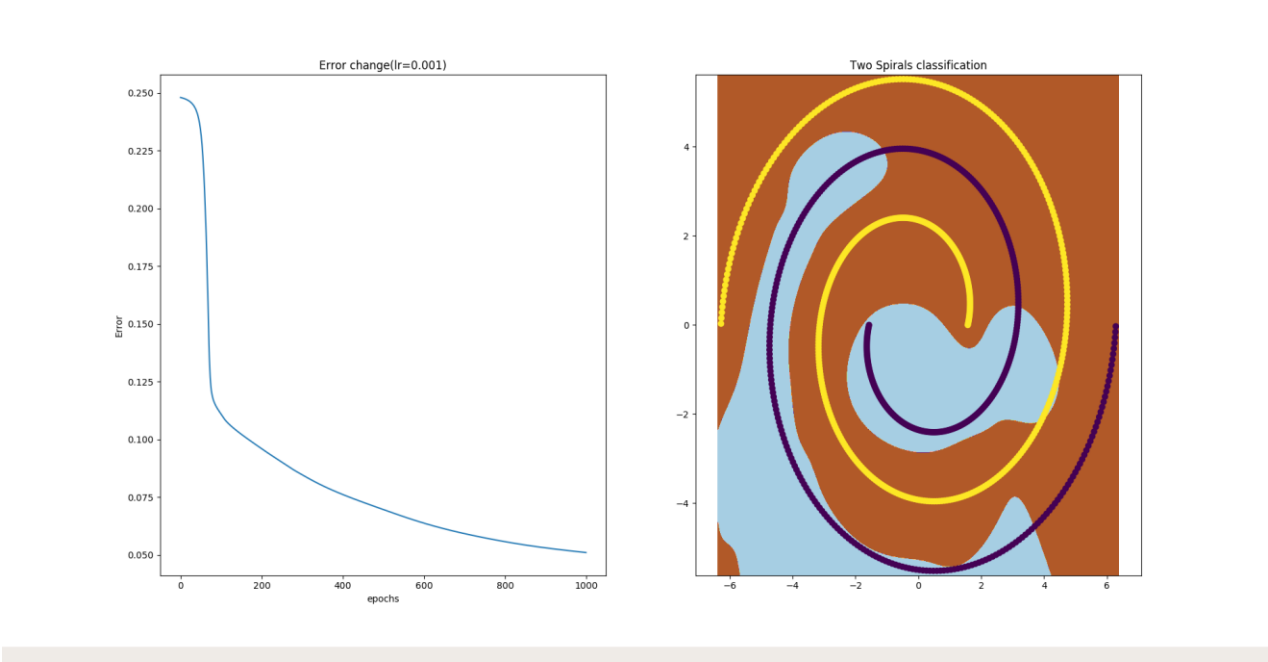
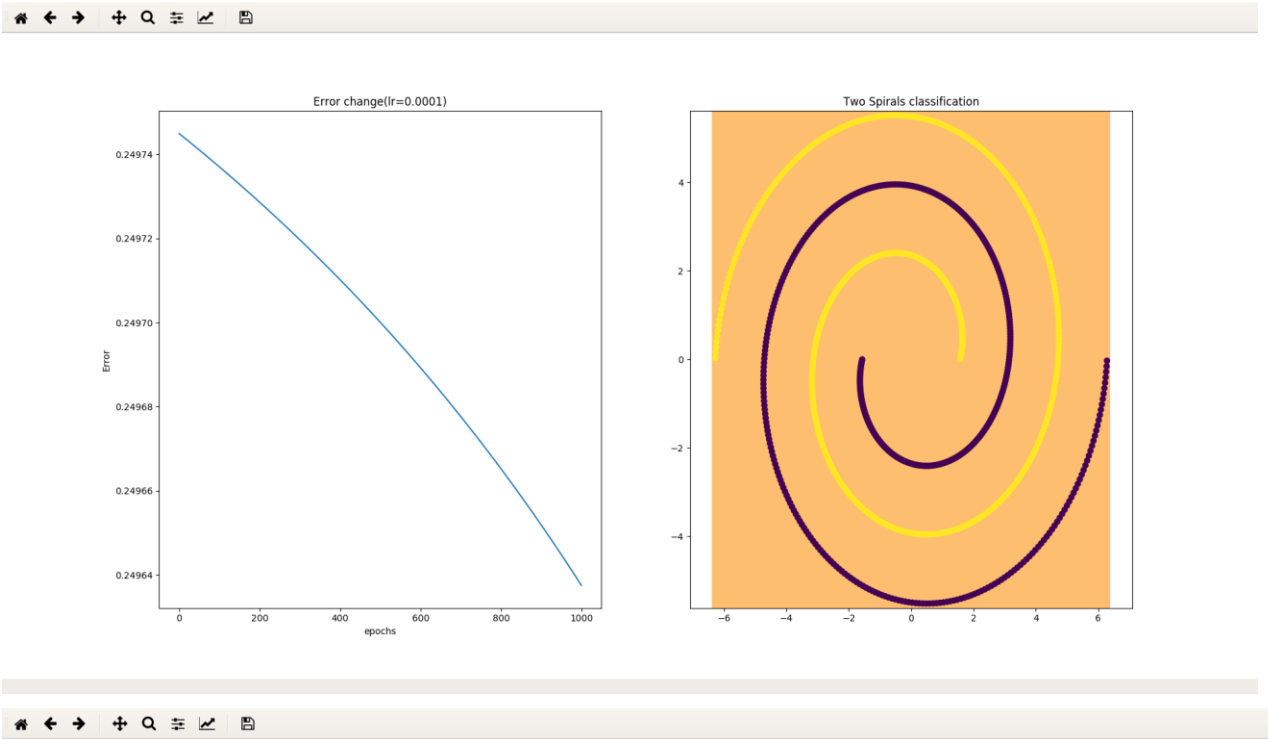
这里我从 csdn 一个聚类数据集里找到了一个双螺旋的数据集，当然也可以根据双螺旋的方程自己 generate 数据点。

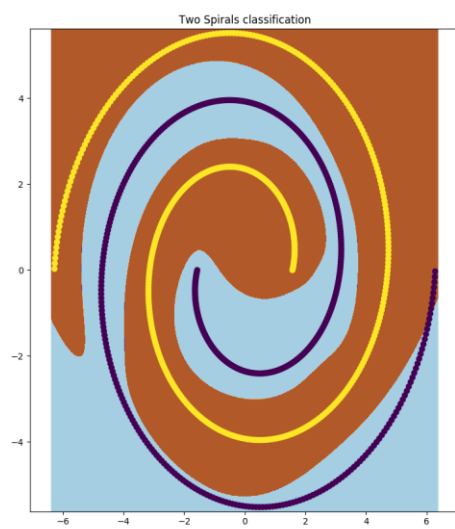
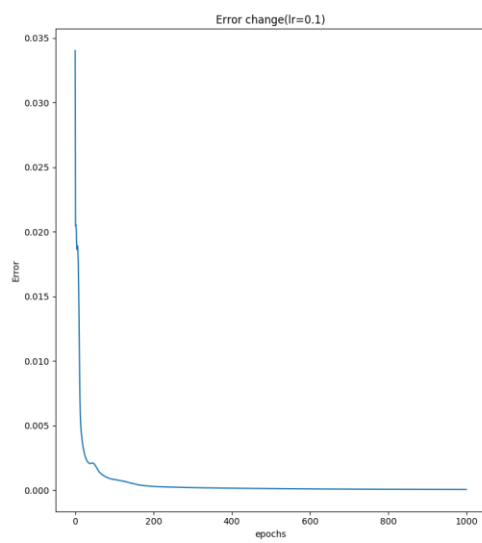
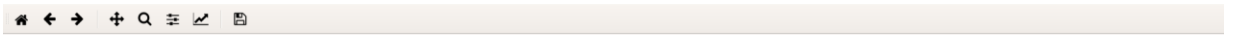
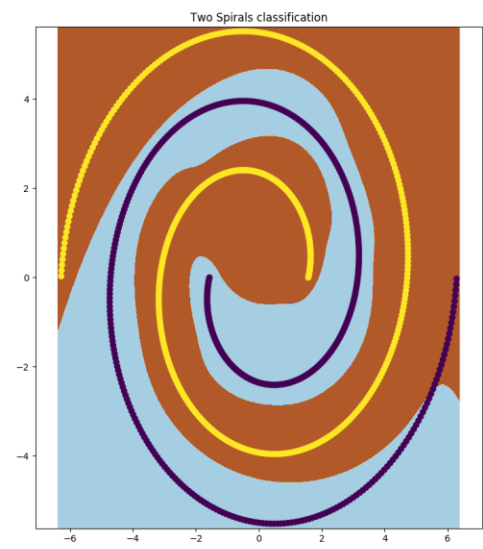
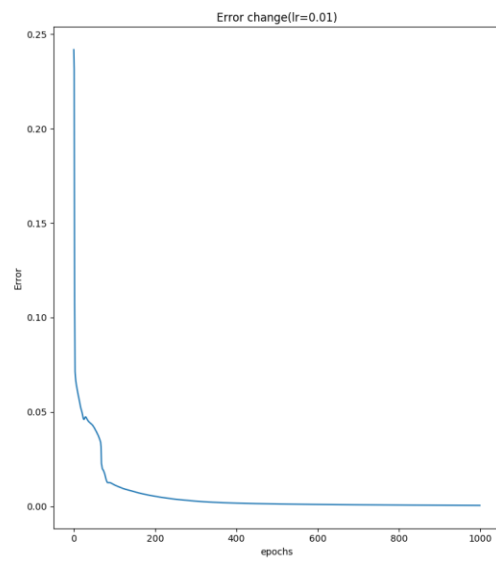
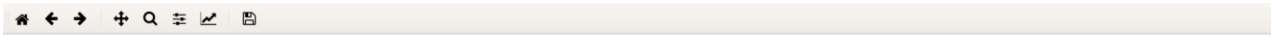
然后构造了神经网络(in:2 hidden:10 out:1)，采用不同的学习率进行训练和测试(learning rate:=0.0001, 0.001, 0.01, 0.1, 0.8)，这里为了更好的比较，我选了 5 个学习率，训练的 epoch 是 1000。

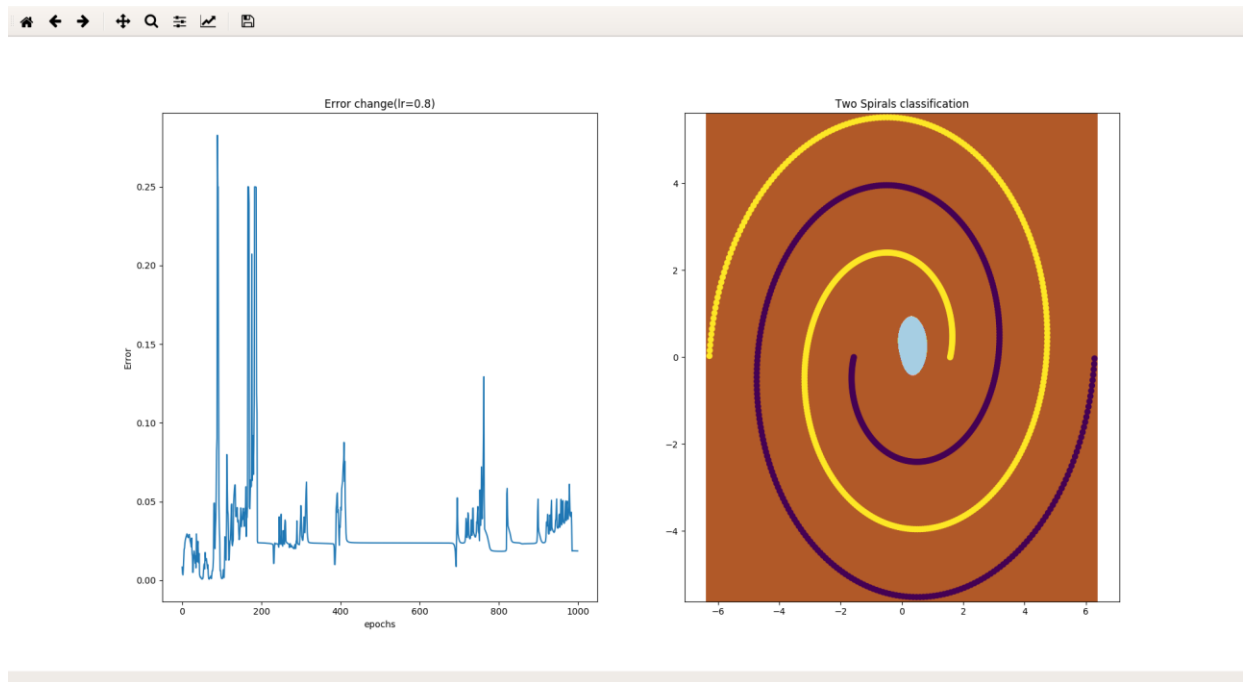
程序主要部分截图如下：

```
13 nn = BP_network() # build a BP network class
14 nn.CreateNN(2, 10, 1, 'Sigmoid') # build the network
15
16 e = []
17 for i in range(1000):
18     print(i)
19     err, err_k = nn.TrainStandard(X, Y.reshape(len(Y), 1), lr=0.01)
20     e.append(err)
21
22 axes[0].set_xlabel("epochs")
23 axes[0].set_ylabel("Error")
24 axes[0].set_title("Error change(lr=0.01)")
25 axes[0].plot(e)
```

训练结果如下（学习率在每张图的左图上方标题标出）：







分析:

首先从训练的时间上来看， $lr=0.0001$ 在 1000epochs 后还没有收敛， $lr=0.001-0.1$ 收敛时间越来越快， $lr=0.8$ 的时候误差抖动很大没有达到收敛。

从决策面的效果来看， $lr=0.0001$ 的决策面还没有训练出来，训练十分缓慢， $lr=0.001$ 决策面稍微有点贴合整体结构了， $lr=0.01$ 基本符合， $lr=0.1$ 训练出来的就很贴合了，然后 $lr=0.8$ 的时候决策面又偏离了数据集的结构特征。

综上 $lr=0.1$ 左右的时候训练效果是比较好的。