

exp 2-分类聚类

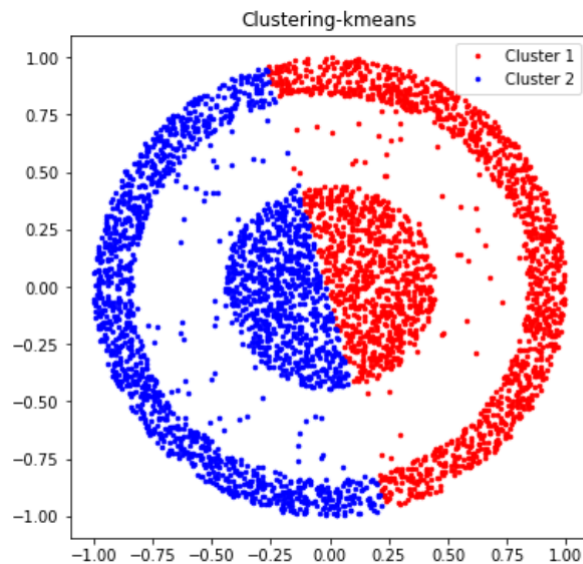
聚类算法实现

k-means算法

实现代码：

```
1  def kmeans(X, k):
2      N, P = X.shape
3      idx = np.zeros(N)
4      # YOUR CODE HERE
5      # -----
6      centers = {}
7      for i in range(k):
8          centers[i] = X[i]
9      # update points
10     for i in range(1000):
11         subsets = {}
12         for j in range(k):
13             subsets[j] = []
14         for point in X:
15             dist = []
16             for center in centers:
17                 dist.append(np.linalg.norm(point-centers[center]))
18     # update centers
19         subsets[dist.index(min(dist))].append(point)
20         pre_centers = dict(centers)
21         for c in subsets:
22             centers[c] = np.average(subsets[c], axis=0)
23     # calculate diff
24         flag = True
25         for c in centers:
26             pre_c = pre_centers[c]
27             cur_c = centers[c]
28             if np.sum((cur_c - pre_c)/pre_c*100) > 0.01:
29                 flag = False
30         if flag:
31             break
32
33     for i in range(N):
34         dist = []
35         for center in centers:
36             dist.append(np.linalg.norm(X[i]-centers[center]))
37         idx[i] = dist.index(min(dist))
38     # -----
39     return idx
```

效果图：



其中起始中心点为文件中从头开始的k个点。

谱聚类算法

实现代码：

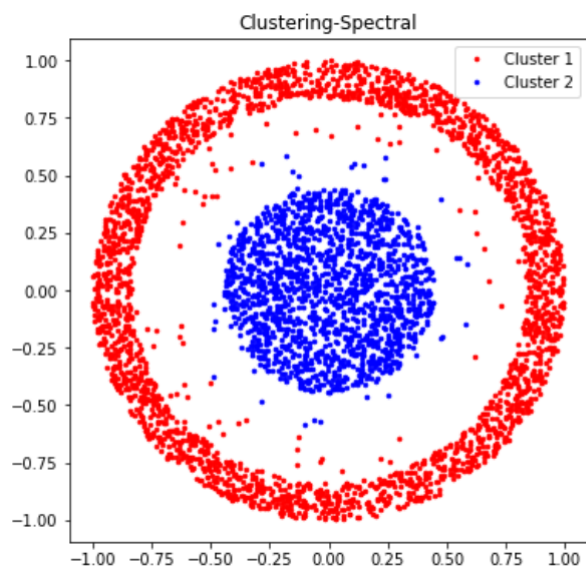
```

1  def spectral(w, k):
2      N = w.shape[0]
3      idx = np.zeros((N, 1))
4      # YOUR CODE HERE
5      # -----
6      # get degree matrix
7      size_ = len(w)
8      D = np.diag(np.zeros(size_))
9      for i in range(size_):
10         D[i][i] = sum(w[i])
11      # get L matrix
12      L = D - w
13      # get eig matrix
14      eigval, eigvec = np.linalg.eig(L)
15      dim = len(eigval)
16      dictEigval = dict(zip(eigval, range(0, dim)))
17      kEig = np.sort(eigval)[0:k]
18      ix = [dictEigval[k] for k in kEig]
19      X = eigvec[:, ix]
20      # -----
21      X = X.astype(float) # keep real part, discard imaginary part
22      idx = kmeans(X, k)
23      return idx

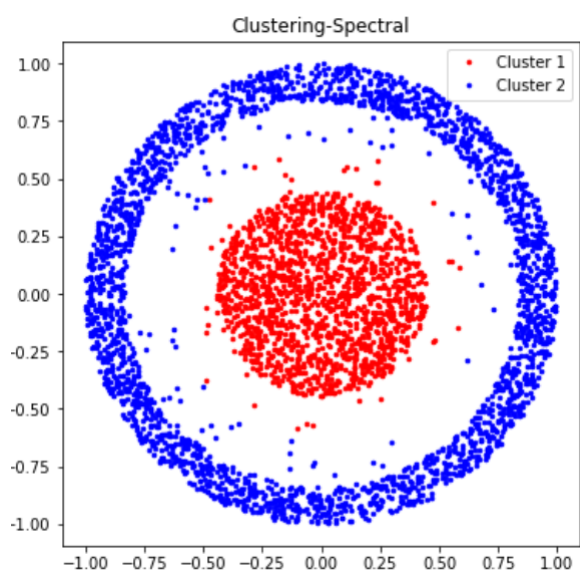
```

效果图：

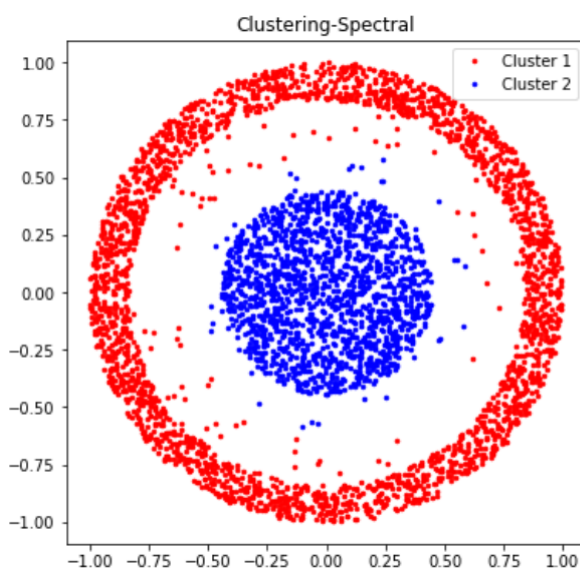
knn参数-k=20, threshold=1.45



knn参数-k=30,threshold=1.45



knn参数-k=10,threshold=1.45



参数threshold在[1,5]区间图像没有明显区别。

两种聚类算法的聚类结果区别

k-means聚类算法只能进行凸函数的聚类，面对像环状的数据集不能很好地区分，谱聚类算法可以对非凸函数进行比较好的聚类。（k-means比较快速，谱聚类耗时较久）

分类算法实现

LR模型

首先，选择划分训练集为0.8*源数据集大小。测试集为剩下的0.2。

先对给定的X进行扩充1使其成为 $[1, x, x^2]$ 的矩阵：

```
1 P, N = X.shape
2 X_t = np.zeros((P+1, N))
3 X_t[0] = 1
4 X_t[1:] = X_train
```

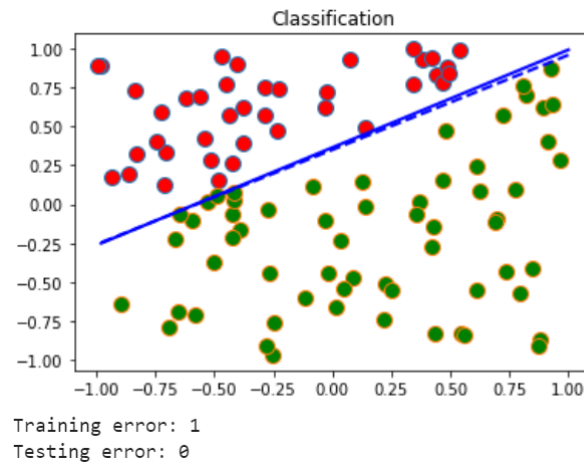
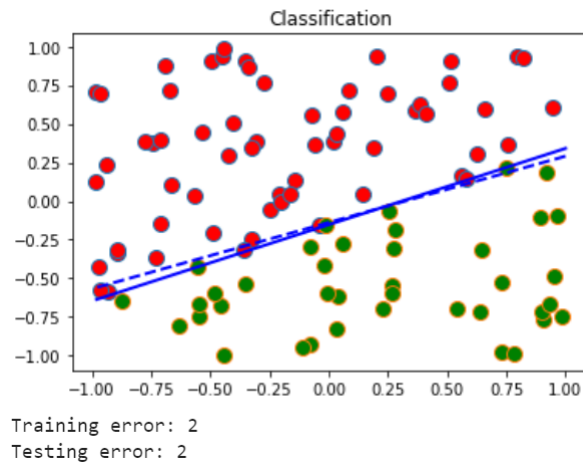
然后使用LR模型公式对权重矩阵进行计算：

```
1 w =
  np.dot(np.dot(np.linalg.inv(np.dot(X_t,X_t.transpose())) ,X_t),y_train.transpo
    se())
```

对训练集进行训练，并对测试集进行测试，记录错误数据：

```
1 train_err = 0
2 test_err = 0
3 P, N = X_train.shape
4 X_t = np.zeros((P+1, N))
5 X_t[0] = 1
6 X_t[1:] = X_train
7 trainy = np.dot(w_1.transpose(),X_t)
8 P, N = X_test.shape
9 X_t = np.zeros((P+1, N))
10 X_t[0] = 1
11 X_t[1:] = X_test
12 testy = np.dot(w_1.transpose(),X_t)
13 for i in range(no_train):
14     if trainy[0][i] * y_train[0][i] < 0:
15         train_err += 1
16 for i in range(no_test):
17     if testy[0][i] * y_test[0][i] < 0:
18         test_err += 1
```

效果如下：



感知机模型

首先，选择划分训练集为0.8*源数据集大小。测试集为剩下的0.2。

先对给定的X进行扩充1使其成为 $[1, x, x^2]$ 的矩阵。

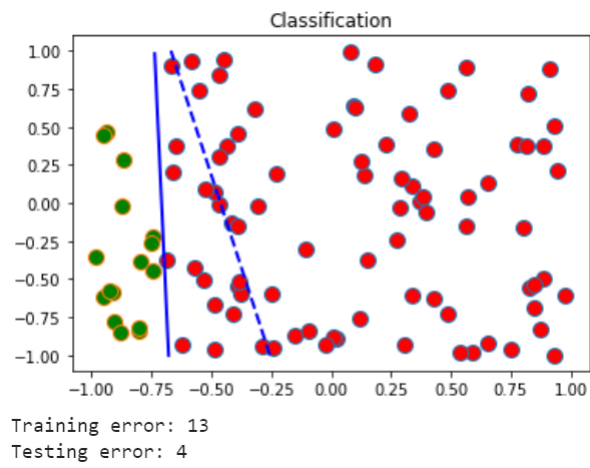
然后对数据集进行训练。这里使用一个神经元，设置学习率为0.75，循环n次，每次对每个输入输出进行预测和调整：

```
1 for i in range(n):
2     for i in range(N):
3         if np.dot(X_t[i], w) * y[i] <= 0:
4             z = learning_rate * X_t[i] * (y[i]-np.dot(X_t[i], w))
5             for j in range(P+1):
6                 w[j] += z[j]
```

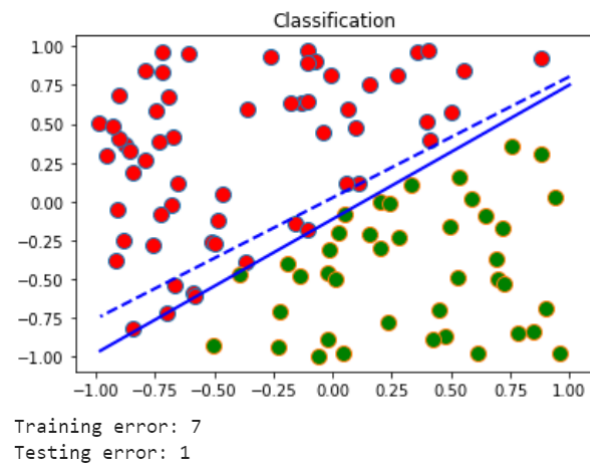
对训练集进行训练，并对测试集进行测试，记录错误数据。

效果如下：

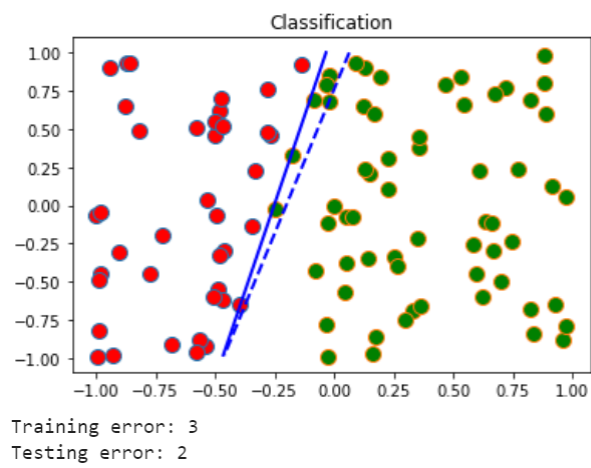
n=1



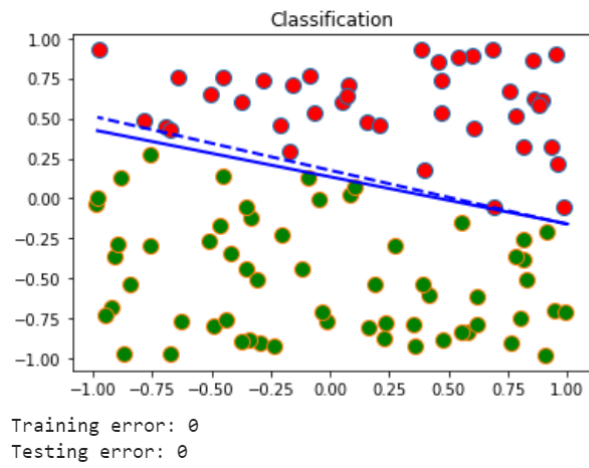
n=10



n=20



n=30



异同分析

LR模型可以通过计算公式得到，所以很快，而感知机在迭代次数少时效果较差，所以需要较多的迭代次数，耗时较长。但是LR模型的准确率没有多次迭代的感知机高

注：代码中默认使用func函数，即LR模型，如果要使用感知机模型，请import func2, 并将func函数名称改为fun2.