

# 模式识别作业8

张蔚桐 2015011493 自55

## 1 MDS算法的应用

采用MDS算法对给出的数据二维化之后可以得到如图1 所示的平面坐标图  
对比图1和实际情况我们可以看出，主要存在几点不同

1. 东西南北位置不对应

这一点是因为MDS算法没有相关的信息而出现的情况

2. 有些城市的相对地理位置不对

这一点可能是因为火车时间间隔和实际的地理情况有关，如西部地区交通比较发达等

## 2 数据降维和处理

我们使用PCA，LLE，TSNE三种方法对问题中给出的数据进行了降维处理并使用相同的（40隐节点的单层神经网络）进行了测试，测试正确率如图2所示，可以看出测试正确率在 98%以上

### 2.1 PCA

方法

PCA 方法的主要流程是将所给数据集进行类似KL变换，得到变换矩阵和对应的方差分量。我们希望保证整个数据集的95%左右的方差，因此PCA采用了前133个变换向量，变换之后的二维展示效果和每个变换对应的比例如图3所示

用和原始同样的神经网络对降维之后的数据加以训练，得到的混淆矩阵如图4所示，可以看出虽然正确率有所下降，但下降并不明显。但特征数是原来的 $\frac{1}{6}$ ，效果还是很满意的。

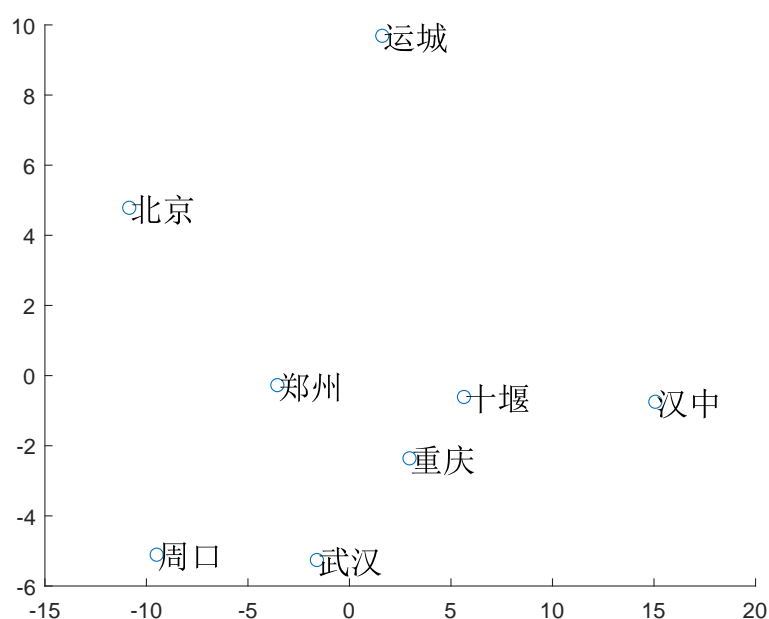


Fig. 1: MDS方法处理过后的地图

## 2.2 LLE方法

LLE方法已经被完整的封装到给定的代码框架内，我们调用之后直接输出前两维的情况即可得到需要的二维展示图，如图5所示

在试图对LLE方法压缩的数据进行训练的过程中我们发现，LLE方法压缩的数据和数据本身有关，因此使用两组分别压缩的测试集和训练集得到的训练——测试的效果很不好，我们对压缩之后的训练集进行分割后测试训练集保留的测试部分正确率如图6所示，可以看出，虽然仅剩下二维特征，但是训练的效果还是比较好的。

## 2.3 tSNE

方法和LLE方法相同，tSNE方法也主要基于给定的算法框架。我们得到二维平面显示如图7所示

和LLE方法相同，这种方法的训练和数据集有关，因此我们需要在训练集上分割出用于测试的部分，得到的训练效果如图8所示，效果还是比较理想的。

总结三种降维方式，PCA线性降维，需要的计算量相对较小，后两种可以处理非线性情况，其中就本题来看tSNE效果最好。

Confusion Matrix						
Output Class	1	2	3	4	5	
	75 15.9%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	0 0.0%	91 19.3%	0 0.0%	1 0.2%	0 0.0%	98.9% 1.1%
	0 0.0%	1 0.2%	90 19.1%	2 0.4%	2 0.4%	94.7% 5.3%
	0 0.0%	0 0.0%	4 0.8%	89 18.9%	0 0.0%	95.7% 4.3%
	0 0.0%	0 0.0%	4 0.8%	0 0.0%	113 23.9%	96.6% 3.4%
Target Class						100% 0.0%
						98.9% 1.1%
						91.8% 8.2%
						96.7% 3.3%
						98.3% 1.7%
						97.0% 3.0%

Fig. 2: 不降维的测试混淆矩阵

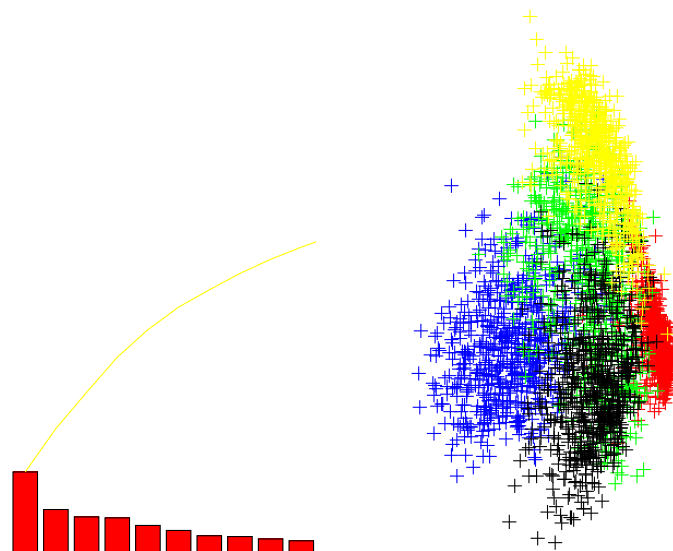


Fig. 3: PCA降维之后的二维平面显示

三种降维均对正确率有损失，但是降维程度较大，小部分的损失是完全可以接受的。

## 2.4 特征选取

本题采用基于类内间隔和类间间隔的选择方法，为加快程序的运行效率，选取了50个特征并采用神经网络加以训练

首先程序将数据集分割为10份，之后每次选取9份进行训练，一份进行测试，训练集上采用 $\text{tr}S_b/\text{tr}S_w$ 的指标进行特征选择并训练网络，十次分别测试得到的特征值。

特征选择过程中，我们认为每个特征之间是独立的，采用类内——类间准则为每个特征评分，选出最合适的前 $n$ 个准则加以组合。

在这种条件下正确率在92%，选取特征保存在index中，发现有很大的集中性，例如5, 48,917号特征常常被选中，说明测试效果良好。

数据保存在ratio.mat 和index.mat中

注：按照要求已经将所有数据集剔除提交，如需重新运行代码需要自行添加数据集，代码运行时间可能较长

Confusion Matrix						
Output Class	1	2	3	4	5	
	74 15.7%	0 0.0%	0 0.0%	1 0.2%	0 0.0%	98.7% 1.3%
	0 0.0%	90 19.1%	1 0.2%	1 0.2%	0 0.0%	97.8% 2.2%
	1 0.2%	1 0.2%	93 19.7%	4 0.8%	2 0.4%	92.1% 7.9%
	0 0.0%	1 0.2%	4 0.8%	86 18.2%	0 0.0%	94.5% 5.5%
	0 0.0%	0 0.0%	0 0.0%	0 0.0%	113 23.9%	100% 0.0%
Target Class						96.6% 3.4%

Fig. 4: PCA方法的测试正确率

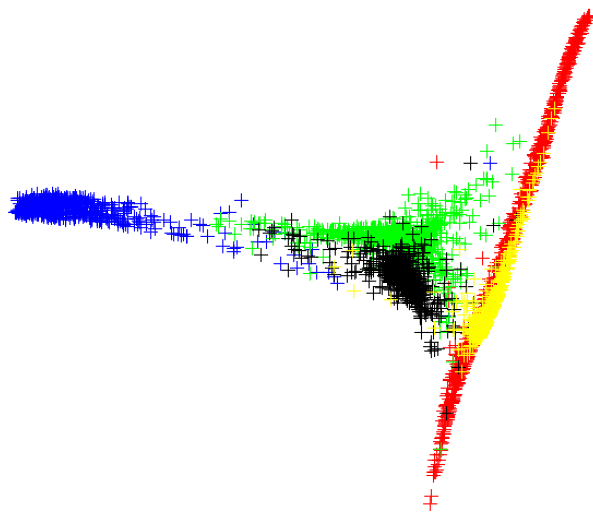


Fig. 5: LLE方法分类图

410	0	4	2	0	98.6%
18.6%	0.0%	0.2%	0.1%	0.0%	1.4%
0	414	5	4	56	86.4%
0.0%	18.8%	0.2%	0.2%	2.5%	13.6%
5	3	378	42	2	87.9%
0.2%	0.1%	17.2%	1.9%	0.1%	12.1%
5	1	31	423	5	91.0%
0.2%	0.0%	1.4%	19.2%	0.2%	9.0%
0	66	1	0	344	83.7%
0.0%	3.0%	0.0%	0.0%	15.6%	16.3%
97.6%	85.5%	90.2%	89.8%	84.5%	89.5%
2.4%	14.5%	9.8%	10.2%	15.5%	10.5%

86	0	0	0	0	100%
18.2%	0.0%	0.0%	0.0%	0.0%	0.0%
0	87	0	3	9	87.9%
0.0%	18.4%	0.0%	0.6%	1.9%	12.1%
4	0	82	11	0	84.5%
0.8%	0.0%	17.4%	2.3%	0.0%	15.5%
1	0	5	80	1	92.0%
0.2%	0.0%	1.1%	16.9%	0.2%	8.0%
0	15	0	2	86	83.5%
0.0%	3.2%	0.0%	0.4%	18.2%	16.5%
94.5%	85.3%	94.3%	83.3%	89.6%	89.2%
5.5%	14.7%	5.7%	16.7%	10.4%	10.8%

90	0	0	0	0	100%
19.1%	0.0%	0.0%	0.0%	0.0%	0.0%
0	98	1	2	10	88.3%
0.0%	20.8%	0.2%	0.4%	2.1%	11.7%
1	2	73	6	0	89.0%
0.2%	0.4%	15.5%	1.3%	0.0%	11.0%
2	0	7	83	1	89.2%
0.4%	0.0%	1.5%	17.6%	0.2%	10.8%
0	13	0	0	83	86.5%
0.0%	2.8%	0.0%	0.0%	17.6%	13.5%
96.8%	86.7%	90.1%	91.2%	88.3%	90.5%
3.2%	13.3%	9.9%	8.8%	11.7%	9.5%

586	0	4	2	0	99.0%
18.6%	0.0%	0.1%	0.1%	0.0%	1.0%
0	599	6	9	75	86.9%
0.0%	19.0%	0.2%	0.3%	2.4%	13.1%
10	5	533	59	2	87.5%
0.3%	0.2%	16.9%	1.9%	0.1%	12.5%
8	1	43	586	7	90.9%
0.3%	0.0%	1.4%	18.6%	0.2%	9.1%
0	94	1	2	513	84.1%
0.0%	3.0%	0.0%	0.1%	16.3%	15.9%
97.0%	85.7%	90.8%	89.1%	85.9%	89.6%
3.0%	14.3%	9.2%	10.9%	14.1%	10.4%

Fig. 6: LLE方法的测试正确率

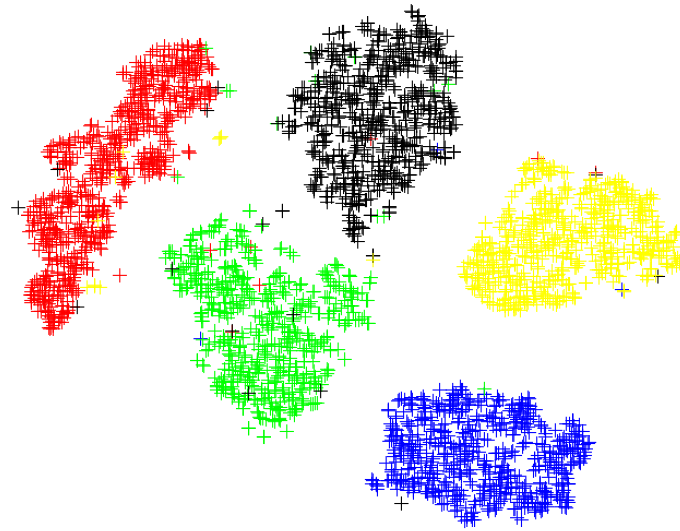


Fig. 7: tSNE方法分类图

417	0	1	1	0	39.5%
18.9%	0.0%	0.0%	0.0%	0.0%	0.0%
0	482	2	4	10	35.8%
0.0%	21.9%	0.1%	0.2%	0.5%	3.2%
1	4	413	6	0	97.4%
0.0%	0.2%	18.8%	0.3%	0.0%	2.6%
1	1	5	444	0	98.4%
0.0%	0.0%	0.2%	20.2%	0.0%	1.6%
1	2	0	2	404	98.8%
0.0%	0.1%	0.0%	0.1%	18.4%	1.2%
98.3%	98.6%	98.1%	97.2%	27.6%	98.1%
0.0%	1.4%	1.8%	2.6%	2.4%	1.9%

88	0	0	0	0	100%
18.6%	0.0%	0.0%	0.0%	0.0%	0.0%
0	113	1	0	0	99.1%
0.0%	23.9%	0.2%	0.0%	0.0%	0.9%
0	0	88	0	0	100%
0.0%	0.0%	18.6%	0.0%	0.0%	0.0%
0	0	3	93	0	98.9%
0.0%	0.0%	0.6%	19.7%	0.0%	3.1%
0	0	0	0	86	100%
0.0%	0.0%	0.0%	0.0%	18.2%	0.0%
100%	100%	95.7%	100%	100%	99.2%
0.0%	0.0%	4.3%	0.0%	0.0%	0.8%

96	0	0	0	0	100%
20.3%	0.0%	0.0%	0.0%	0.0%	0.0%
0	97	1	1	0	98.0%
0.0%	33.6%	0.2%	0.2%	0.0%	2.0%
0	0	73	2	1	98.1%
0.0%	0.0%	15.5%	0.4%	0.2%	3.9%
0	0	0	105	0	100%
0.0%	0.0%	0.0%	22.2%	0.0%	0.0%
0	0	0	0	96	100%
0.0%	0.0%	0.0%	0.0%	20.3%	0.0%
100%	100%	98.6%	97.2%	29.0%	98.9%
0.0%	0.0%	1.4%	2.8%	1.0%	1.1%

601	0	1	1	0	99.7%
19.1%	0.0%	0.0%	0.0%	0.0%	0.3%
0	692	4	5	10	97.3%
0.0%	22.0%	0.1%	0.2%	0.3%	2.7%
1	4	574	8	1	97.6%
0.0%	0.1%	18.3%	0.3%	0.0%	2.4%
1	1	8	642	0	98.5%
0.0%	0.0%	0.3%	20.4%	0.0%	1.5%
1	2	0	2	596	99.2%
0.0%	0.1%	0.0%	0.1%	19.6%	0.8%
99.5%	99.0%	97.6%	97.6%	88.2%	98.4%
0.5%	1.0%	2.2%	2.4%	1.8%	1.6%

Fig. 8: tSNE方法的测试正确率