

# PCA特征降维&J48特征选择算法实现报告

周政

---

## 1.PCA算法

协方差矩阵及其特征向量和特征值计算函数

```
#协方差矩阵及其特征向量和特征值计算函数
def CovFeatureVecsVals(X):
    mean_vec = np.mean(X, axis=0)
    #归一化
    for i in range(len(X)):
        X[i,:] = (X[i,:] - mean_vec[:])
    cov_mat = np.cov(X.T)
    eig_vals, eig_vecs = np.linalg.eig(cov_mat)
    #组成特征值和特征向量对(通过numpy.linalg.eig()计算出的特征值和特征向量是从大到小对应排好序的.)
    eig_pairs = [(eig_vals[i], eig_vecs[:,i]) for i in range(len(eig_vals))]
    return eig_pairs
```

传入数据集X，将特征列减去其均值减少差异程度，然后计算各特征的协方差矩阵及其特征向量与特征值，组成协方差矩阵的特征值与特征向量的tuple组返回

降维矩阵计算函数

```
15 #降维矩阵计算函数:
16 def ReduceDimension(eig_pairs,d):
17     matrix_w = np.hstack(
18         (eig_pairs[i][1].reshape(4,1) for i in range(d))
19     )
20     return matrix_w
```

传入装有特征值与特征向量对的list,以及要降的维度d,来构建降维矩阵并返回，得到的降维矩阵是一个nxd的矩阵。

PCA降维调度函数

```
21 #PCA降维调度函数
22 def RunPCA(data,d):
23     X = np.array(data,dtype = 'float64')
24     d = int(d)
25     if d >= X.shape[1] or d <= 0:
26         return 0
27     #计算协方差矩阵及其特征向量和特征值
28     eig_pairs = CovFeatureVecsVals(X)
29     #计算降维矩阵
30     matrix_w = ReduceDimension(eig_pairs,d)
31     #降维
32     Y = X.dot(matrix_w)
33     return Y
34
```

通过传入初始数据集和需要降维的维数对各函数进行调用最后计算出降维后的矩阵Y(mxd).

Python文件执行程序及画图

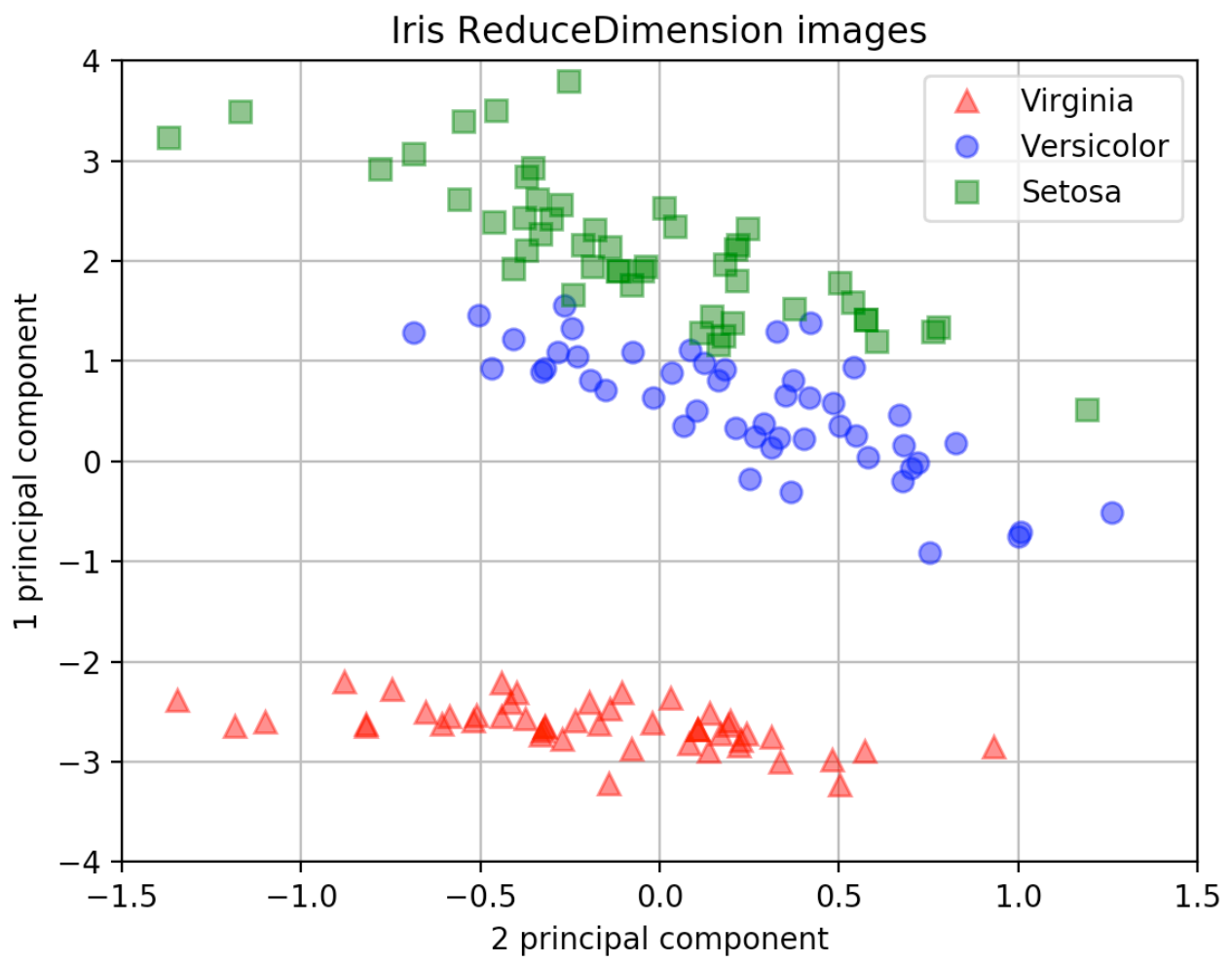
```

36 if __name__ == '__main__':
37     #导入数据
38     data = []
39     with open('dataset/iris.data') as f:
40         for i in f.readlines():
41             data.append(i.strip().split(',')[::-1])
42     #运行PCA调度函数
43     Y = RunPCA(data,2)
44     if type(Y) == type(0):
45         print("无法正确降维! ")
46     else:
47         print(Y)
48         plt.plot(Y[0:49,1],Y[0:49,0], '^', markersize=7, color='red', alpha=0.5,
49                 label='Virginia')
50         plt.plot(Y[50:99,1],Y[50:99,0], 'o', markersize=7, color='blue', alpha=0.5,
51                 label='Versicolor')
52         plt.plot(Y[100:149,1],Y[100:149,0], 's', markersize=7, color='green', alpha=0.5,
53                 label='Setosa')
54
55         plt.xlim([-1.5,1.5])
56         plt.ylim([-4,4])
57
58         plt.xlabel('2 principal component')
59         plt.ylabel('1 principal component')
60         plt.legend()
61         plt.title('Dois componentes principais da Base de Dados Iris')
62
63         plt.tight_layout
64         plt.grid()
65         plt.show()

```

结果

Figure 1



x=-0.335685 y=-1.76623

## 1.利用J48进行特征选择

之前作业写好的J48决策树

#香农熵计算

```
def ShannonEntropy(dataSet):  
    labelCounts = {} #记录标签出现次数  
    dataNum = len(dataSet) #数据条数  
    for data in dataSet:  
        cLabel = data[-1]  
        if cLabel not in labelCounts.keys():  
            labelCounts[cLabel] = 0  
        labelCounts[cLabel] += 1  
    SE = 0.0  
    for key in labelCounts:  
        percent = float(labelCounts[key])/dataNum  
        SE -= percent*log(percent,2)  
    return SE
```

#划分数数据集

```
def DiviDataSet(dataSet, feature, v):  
    dividDataSet = []  
    for data in dataSet:  
        if data[feature] == v:  
            tempDataSet = data[:feature]  
            tempDataSet.extend(data[feature+1:])  
            dividDataSet.append(tempDataSet)  
    return dividDataSet
```

```
26 #选择根节点属性来划分数数据集
27 def selectBestFeature(dataSet):
28     numAttr = len(dataSet[0]) - 1
29     baseEnt = ShannonEntropy(dataSet)
30     bestInfoInc = 0.0
31     bestFeature = -1
32     for i in range(numAttr):
33         uValues = set([example[i] for example in dataSet])
34         newEnt = 0.0
35         for value in uValues:
36             subDataSet = DiviDataSet(dataSet, i, value)
37             percent = len(subDataSet)/float(len(dataSet))
38             newEnt += percent * ShannonEntropy(subDataSet)
39             InfoInc = baseEnt - newEnt
40             if (InfoInc > bestInfoInc):
41                 bestInfoInc = InfoInc
42                 bestFeature = i
43     return bestFeature
```

```
44 #计算出现次数最多的分类
45 def majorityCnt(classList):
46     classCount = {}
47     for v in classList:
48         if v not in classCount.keys():
49             classCount[v] = 0
50             classCount[v] += 1
51     sortedClassCount = sorted(classCount.iteritems(), key=operator.itemgetter(1),
52                               reverse=True)
53     return sortedClassCount[0][0]
```

```

53 #建树
54 def createTree(dataSet, labels, n):
55     classList = [example[-1] for example in dataSet]
56     if classList.count(classList[0]) == len(classList):
57         return classList[0]
58     if len(dataSet[0]) == 1:
59         return majorityCnt(classList)
60     bestFeat = selectBestFeature(dataSet)
61     bestFeatLabel = labels[bestFeat]
62     if bestFeatLabel not in FeatureSubset:
63         if len(FeatureSubset) == n:
64             return
65         FeatureSubset.append(bestFeatLabel)
66     DTree = {bestFeatLabel: {}}
67     del(labels[bestFeat])
68     featValues = [example[bestFeat] for example in dataSet]
69     uniqueVals = set(featValues)
70     for value in uniqueVals:
71         subLabels = labels[:]
72         DTree[bestFeatLabel][value] = createTree(DiviDataSet(dataSet, bestFeat,
73             value), subLabels, n)
73     return DTree

```

但对原来的递归建树函数做了一点修改，新加了传入参数n来获取前n层节点所用到的特征装入FeatureSubset中并在FeatureSubset长度到达n后return不再递归返回结果。

结果

```

Last login: Tue Aug 22 12:20:13 on ttys003
Zain@MacBook-Pro ~ wb/J48\&PCA
Zain@MacBook-Pro ~/wb/J48&PCA  master  py3 J48FeatureSelector.py
['safety', 'persons', 'buying', 'maint']
Zain@MacBook-Pro ~/wb/J48&PCA  master  _

```

[Github地址](#)