

## 目录

1、概念理解.....	3
1.1、机器学习基本概念理解.....	3
1.2、数据分析概念理解.....	4
2、算法理解.....	5
2.1、聚类方法.....	5
2.1.1、K-means.....	5
2.2、分类器.....	6
2.2.1、C4.5 决策树.....	7
2.2.2、KNN.....	8
2.2.3、朴素贝叶斯.....	9
2.3、统计学习.....	9
2.3.1、SVM.....	9
2.4、集成学习.....	9
2.4.1、AdaBoost.....	10
2.5、线性模型.....	11
3、算法实现.....	11
3.1、KNN 算法.....	11
3.1.1、代码实现：.....	11
数据集一：Iris.....	14
数据集二：letter-recognition.....	15
数据集三：cmc.....	16
数据集四：Winequality-red.....	18
3.1.2、小结.....	18
3.2、C4.5 决策树.....	19
3.2.1、代码实现.....	19
数据集一：Car Evaluation.....	25
数据集二：letter-recognition.....	26
数据集三：winequality-red.....	27
数据集四：Fertility.....	28
3.2.2、小结.....	29
3.3、聚类算法 k-means.....	30
3.3.1、代码实现.....	30

数据集一：Iris.....	34
数据集二：winequality-red .....	35
数据集三：seeds.....	35
3.3.2、小结.....	36
3.4、朴素贝叶斯分类器 .....	37
3.4.1、代码实现.....	37
数据集一：Phishing.....	40
3.4.2、小结.....	41
4、总结 .....	41

# 1、概念理解

## 1.1、机器学习基本概念理解

编号	色泽	根蒂	含糖率	好瓜
1	青绿	蜷缩	0.46	是
2	乌黑	蜷缩	0.376	是
3	青绿	硬挺	0.268	否
4	乌黑	稍卷	0.091	否

**数据集：**这组记录的集合称为一个“数据集”。

**样本(sample)：**其中每条记录是关于一个事件或对象（这里是一个西瓜）的描述，称为一个“示例”(instance)或“样本”(sample)。

**标签：**又称标记，是训练样本的“结果”，如上表的“好瓜”。

**属性(attribute)：**反映事件或对象在某方面的表现或性质的事项，称为属性(attribute)或特征(feature)。属性上的取值称为“属性值”(attribute value)。

**连续属性：**属性值为一系列连续的数值，如上表含糖率。

**离散属性：**属性值为离散值，如上表色泽，根蒂。

**属性空间：**由属性张成的空间称为“属性空间”(attribute space)、“样本空间”(sample space)或“输入空间”。

**特征向量：**样本在属性空间中所对应的向量。

**训练集：**从数据中学得模型的训练过程所使用的数据的集合。

**训练样本：**训练集中的样本。

**测试集：**使用模型进行预测时被预测的样本的集合。

**分类：**预测过程所有预测的标记是离散值的学习任务称为分类。如上表西瓜的分类。

**二分类：**只有两个类别，如上表的西瓜分类。两个类分别称为正类和反类。

**聚类：**将训练集的西瓜分为若干组，每一组称为一个“簇”(cluster)，这些自动形成的簇可能对应一些潜在的概念划分，如“浅色瓜”，“深色瓜”，学习过程中所使用的样本通常不含标记信息。

**回归：**预测的是连续的值，可以用直线或者曲线来拟合。如，西瓜的成熟度。

**泛化能力：**即抵抗过拟合的能力，学得模型适用于新样本的能力。

**监督学习：**即训练时，训练样本中含有标签，知道分类的结果。通过答案反向解题。

**非监督学习：**训练样本中不含有标签，不知道分类的结果。

**标准化：**将数据比例缩放，使之落入一个小的待定区域，一般是 $[-1,1]$ ，常用的标准化方法有 min-max 和 z-score。

## 1.2、数据分析概念理解

数据分析是统计学的产物，顾名思义，就是分析数据。数据分析可帮助人们作出判断，以便采取适当行动。在如今这个大数据环境下，对数据分析靠人力是不现实的，使用计算机通过机器学习我想是一个分析数据的很好的办法。

## 2、算法理解

### 2.1、聚类方法

聚类分析以相似性为基础，在一个聚类中的模式之间比不在同一聚类中的模式之间具有更多的相似性。

#### 2.1.1、K-means

描述： k-means 是划分方法中较经典的聚类算法之一。由于该算法的效率高，所以在对大规模数据进行聚类时被广泛应用。目前，许多算法均围绕着该算法进行扩展和改进。

k-means 算法目标是，以 k 为参数，把 n 个对象分成 k 个簇，使簇内具有较高的相似度，而簇间的相似度较低。

k-means 算法的处理过程如下：首先，随机地 选择 k 个对象，每个对象初始地代表了一个簇的平均值或中心;对剩余的每个对象，根据其与各簇中心的距离，将它赋给最近的簇;然后重新计算每个簇的平均值。这个过程不断重复，直到准则函数收敛。通常，采用平方误差准则，其定义如下：

$$E(E = \sum_{i=1}^k \sum_{x \in C_i} \|x - u_i\|^2)$$

这里 E 是数据库中所有对象的平方误差的总和，x 是空间中的点， $u_i$  是簇  $C_i$  的平均值。该目标函数使生成的簇尽可能紧凑独立，使用的距离度量是欧几里得距离,当然也可以用其他距离度量。k-means 聚类算法的算法流程如下：

输入：包含 n 个对象的数据库和簇的数目 k;

输出：k 个簇，使平方误差准则最小。

步骤:

- (1) 任意选择  $k$  个对象作为初始的簇中心;
- (2) repeat;
- (3) 根据簇中对象的平均值, 将每个对象(重新)赋予最类似的簇;
- (4) 更新簇的平均值, 即计算每个簇中对象的平均值;
- (5) until 不再发生变化。

总结:

优点: 简单直接 (体现在逻辑思路以及实现难度上), 易于理解, 在低维数据集上有不错的效果。

缺点: 对于高维数据 (如成百上千维, 现实中还不止这么多), 其计算速度十分慢, 主要是慢在计算距离上 (参考欧几里得距离, 当然并行化处理是可以的, 这是算法实现层面的问题), 它的另外一个缺点就是它需要我们设定希望得到的聚类数  $k$ , 若我们对于数据没有很好的理解, 那么设置  $k$  值就成了一种估计性的工作。

## 2.2、分类器

分类是数据挖掘的一种非常重要的方法。分类的概念是在已有数据的基础上学会一个分类函数或构造出一个分类模型 (即我们通常所说的分类器 (Classifier))。该函数或模型能够把数据库中的数据纪录映射到给定类别中的某一个, 从而可以应用于数据预测。总之, 分类器是数据挖掘中对样本进行分类的方法的统称, 包含决策树、逻辑回归、朴素贝叶斯、神经网络等算法。

### 2.2.1、C4.5 决策树

**简介：**C4.5 算法是用于生成决策树的一种经典算法，是 ID3 算法的一种延伸和优化。C4.5 算法对 ID3 算法主要做了一下几点改进：

- (1) 通过信息增益率选择分裂属性，克服了 ID3 算法中通过信息增益倾向于选择拥有多个属性值的属性作为分裂属性的不足；
- (2) 能够处理离散型和连续型的属性类型，即将连续型的属性进行离散化处理；
- (3) 避免过拟合，构造决策树之后进行剪枝操作；
- (4) 能够处理具有缺失属性值的训练数据。

**分裂属性选择：**分裂属性选择的评判标准是决策树算法之间的根本区别。区别于 ID3 算法通过信息增益选择分裂属性，C4.5 算法通过信息增益率选择分裂属性。

**信息熵：**

是度量样本集合纯度最常用的一种指标，假定当前样本集合  $D$  中第  $k$  类样本所占比例为  $p_k$  ( $k = 1, 2, \dots, |Y|$ )，则  $D$  的信息熵定义为：

$$Ent(D) = \sum_{k=1}^{|Y|} p_k \log_2 p_k$$

**信息增益：**

$$Gain(D, a) = Ent(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} Ent(D)^v$$

假定离散属性  $a$  有  $V$  个可能的取值，若使用  $a$  来对样本集  $D$  进行划分，则会产生  $V$  个分支结点，其中第  $v$  个分支结点包含了  $D$  中所有在属性  $a$  上取值为

$a^v$  的样本，记为  $D^v$ 。

信息增益率准则对可取值数目较多的属性有所偏好，所以提出增益率(gain ratio)

$$Gain\_ratio(D, a) = \frac{Gain(D, a)}{IV(a)}$$

其中

$$IV(a) = - \sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$$

称为属性  $a$  的“固有值”。属性  $a$  的可能性取值数目越多（即  $V$  越大），则  $IV(a)$  的值通常会增大。但增益率准则对可取值数目较少的属性有所偏好，因此，C4.5 算法使用启发式：先从候选划分属性中找出信息增益高于平均水平的属性，再从中选择增益率最高的。

Tips：连续属性与离散属性的信息增益计算公式有所不同。

## 2.2.2、KNN

KNN 算法又称  $k$  近邻分类(k-nearest neighbor classification)算法。它是根据不同特征值之间的距离来进行分类的一种简单的机器学习方法，它是一种简单但是懒惰的算法。他的训练数据都是有标签的数据，即训练的数据都有自己的类别。KNN 算法主要应用领域是对未知事物进行分类，即判断未知事物属于哪一类，判断思想是，基于欧几里得定理，判断未知事物的特征和哪一类已知事物的特征最接近。它也可以用于回归，通过找出一个样本的  $k$  个最近邻居，将这些邻居的属性的平均值赋给该样本，就可以得到该样本的属性。K 的取值直接影响最后的结果。



### 2.2.3、朴素贝叶斯

基于概率论的分类算法，原理简单，多用于文本分类。利用贝叶斯公式：

$$P(B_i|A) = \frac{P(B_i)P(A|B_i)}{\sum_{j=1}^n P(B_j)P(A|B_j)}$$

**假设基础：**对于每一个特征都有独立、相等，导致难以找到合适的数据集。

**缺点：**需要知道先验概率，且先验概率很多时候取决于假设，假设的模型可以有很多种，因此在某些时候会由于假设的先验模型的原因导致预测效果不佳。

## 2.3、统计学习

### 2.3.1、SVM

SVM 的分类思想本质上和线性回归 LR 分类方法类似，就是求出一组权重系数，在线性表示之后可以分类。先使用一组 training set 来训练 SVM 中的权重系数，然后可以对 testingset 进行分类。

## 2.4、集成学习

集成学习(ensemble learning)，它本身不是一个单独的机器学习算法，而是通过构建并结合多个机器学习器来完成学习任务。集成学习有两个主要的问题需要解决，第一是如何得到若干个个体学习器，第二是如何选择一种结合策略，将这些个体学习器集成成一个强学习器。

个体学习器有两种，第一种就是所有的个体学习器都是一个种类的，或者说是同质的。

第二种是所有的个体学习器不全是一个种类的，或者说是异质的。

同质个体学习器按照个体学习器之间是否存在依赖关系可以分为两类：

第一个是个体学习器之间存在强依赖关系，一系列个体学习器基本都需要串行生成，代表算法是 boosting 系列算法。

第二个是个体学习器之间不存在强依赖关系，一系列个体学习器可以并行生成，代表算法是 bagging 和随机森林 (Random Forest) 系列算法。

### 2.4.1、AdaBoost

Boosting, 也称为增强学习或提升法, 是一种重要的集成学习技术, 能够将预测精度仅比随机猜度略高的弱学习器增强为预测精度高的强学习器, 这在直接构造强学习器非常困难的情况下, 为学习算法的设计提供了一种有效的新思路和新方法。其中最为成功应用的是, Yoav Freund 和 Robert Schapire 在 1995 年提出的 AdaBoost 算法。

AdaBoost 是英文 "Adaptive Boosting" (自适应增强) 的缩写, 它的自适应在于: 前一个基本分类器被错误分类的样本的权值会增大, 而正确分类的样本的权值会减小, 并再次用来训练下一个基本分类器。同时, 在每一轮迭代中, 加入一个新的弱分类器, 直到达到某个预定的足够小的错误率或达到预先指定的最大迭代次数才确定最终的强分类器。

Adaboost 算法可以简述为三个步骤:

- (1) 首先, 是初始化训练数据的权值分布  $D_1$ 。假设有  $N$  个训练样本数据, 则每一个训练样本最开始时, 都被赋予相同的权值:  $w_1 = 1/N$ 。
- (2) 然后, 训练弱分类器  $h_i$ 。具体训练过程中是: 如果某个训练样本点, 被弱分类器  $h_i$  准确地分类, 那么在构造下一个训练集中, 它对应的权值要减小; 相

反，如果某个训练样本点被错误分类，那么它的权值就应该增大。权值更新过的样本集被用于训练下一个分类器，整个训练过程如此迭代地进行下去。

(3) 最后，将各个训练得到的弱分类器组合成一个强分类器。各个弱分类器的训练过程结束后，加大分类误差率小的弱分类器的权重，使其在最终的分类函数中起着较大的决定作用，而降低分类误差率大的弱分类器的权重，使其在最终的分类函数中起着较小的决定作用。

也就是说，误差率低的弱分类器在最终分类器中占的权重较大，否则较小。

## 2.5、线性模型

线性回归是一种监督学习下的线性模型，线性回归试图从给定数据集中学习一个线性模型来较好的预测输出。其中只需要学得  $w$  和  $b$  就可以了。使用最小二乘法求解  $w$  和  $b$ ，即基于均方误差最小化来进行模型求解。

# 3、算法实现

**说明：**代码实现均使用的 UCI 上数据集。采用  $k$  折交叉验证，验证自己算法模型的稳定性。由于  $K$  折交叉验证是最后才使用的，所以 PDF 里的截图可能没有更新，不过在具体的代码文件里均已更新。

## 3.1、KNN 算法

### 3.1.1、代码实现：

**分类：**

```

1  import numpy as np
2  import random
3
4
5  def deal_with_res(k):
6      lines_set = open('E:\myCodes\python\mlTest\venv\src\d'
7                      'ata\iris.data.txt').readlines()
8      data = []
9      for temp in lines_set:
10         temp2 = list(temp[:].split("\n")[0].split(","))
11         tmp = []
12         for i in range(len(temp2)):
13             if i != len(temp2) - 1:
14                 tmp.append(float(temp2[i]))
15             else:
16                 tmp.append(temp2[i])
17         data.append(tmp[:])
18     list1 = list(range(0, len(data), k))
19     w = random.sample(list1, 1)[0]
20     data_x = []
21     data_x.extend(data[:w])
22     data_x.extend(data[w + int(len(data) / k):])
23     data_y = [i[-1:] for i in data_x]
24     data_x = [i[:-1] for i in data_x]
25     test_y = [i[-1:] for i in data[w:w + int(len(data) / k)]]
26     test_x = [i[:-1] for i in data[w:w + int(len(data) / k)]]
27     return data_x, data_y, test_x, test_y
28

```

```

25     temp = [list_sub(temp, inputs) for temp in test_x]
26     temp2 = np.sum(temp, axis=1)
27     distance = np.ndarray.tolist(temp2 ** 0.5)
28     # 返回排序后的索引值
29     sorted_dist_indices = arg_sort(distance)
30     class_count = {}
31     # 选择k个最近邻
32     for k in range(k):
33         vote_labels = test_y[sorted_dist_indices[k][0]]
34         temp = str(vote_labels)[2:-2]
35         class_count[temp] = class_count.get(temp, 0) + 1
36     max_count = 0
37     # 找出出现次数最多的类别

```

```

38     for k in class_count:
39         if class_count.get(k) > max_count:
40             max_count = class_count.get(k)
41             max_key = k
42     return max_key
43
44
45 def list_sub(a, b):
46     a = list(map(float, a))
47     b = list(map(float, b))
48     return list(np.fabs(np.array(a) - np.array(b)))
49
50
51 def calc(test_x, test_y, inputs, k):
52     # 计算距离
53     temp = [list_sub(temp, inputs) for temp in test_x]
54     temp2 = np.sum(temp, axis=1)
55     distance = np.ndarray.tolist(temp2 ** 0.5)
56     # 返回排序后的索引值
57     sorted_dist_indices = arg_sort(distance)
58     class_count = {}
59     # 选择k个最近邻
60     for a in range(k):
61
62         vote_labels = test_y[sorted_dist_indices[a][0]]
63         temp = str(vote_labels)[2:-2]
64         class_count[temp] = class_count.get(temp, 0) + 1
65
66     max_count = 0
67     for c in class_count:
68         if class_count.get(c) > max_count:
69             max_count = class_count.get(c)
70             max_key = c
71     return max_key
72
73
74 def arg_sort(distance):
75     return sorted(enumerate(distance), key=lambda x: x[1])
76

```

```

87 ▶ if __name__ == '__main__':
88     # list1[0] list1[1] list1[2] list1[3]
89     # 分别代表样本属性，样本标签，测试属性，测试标签
90     p = 10
91     K = 5
92     avg_accuracy = 0
93     for w in range(p):
94         accuracy = 0
95         list1 = deal_with_res(p)
96         total = len(list1[2])
97         count = 0
98         for i in range(total):
99             prdict = calc(list1[0], list1[1], list1[2][i], K)
100             if prdict == "".join(list1[3][i]):
101                 count = count + 1
102             lebal_count = set(j[0] for j in list1[1])
103             accuracy = count / total * 100
104             avg_accuracy += accuracy / p
105             print("第", (w + 1), "次测试准确率为: ", accuracy, "%")
106             if w == p-1:
107                 print("样本数量为: ", len(list1[0]))
108                 print("测试数量为: ", len(list1[2]))
109                 print("属性类为: ", len(list1[0][0]))
110                 print("标签结果: ", len(lebal_count))
111                 print("K值为: ", K)
112                 print("测试平均准确率为: ", avg_accuracy, "%")
113

```

## 数据集一：Iris

```

1      5.0, 3.6, 1.4, 0.2, Iris-setosa
2      5.4, 3.9, 1.7, 0.4, Iris-setosa
3      4.7, 3.2, 1.3, 0.2, Iris-setosa
4      4.6, 3.1, 1.5, 0.2, Iris-setosa
5      6.6, 2.9, 4.6, 1.3, Iris-versicolor
6      5.2, 2.7, 3.9, 1.4, Iris-versicolor
7      6.1, 2.9, 4.7, 1.4, Iris-versicolor
8      4.4, 2.9, 1.4, 0.2, Iris-setosa
9      6.4, 3.2, 5.3, 2.3, Iris-virginica
10     6.5, 3.0, 5.5, 1.8, Iris-virginica
11     6.4, 3.1, 5.5, 1.8, Iris-virginica

```

## 结果:

```
第 1 次测试准确率为: 86.66666666666667 %
第 2 次测试准确率为: 100.0 %
第 3 次测试准确率为: 100.0 %
第 4 次测试准确率为: 86.66666666666667 %
第 5 次测试准确率为: 100.0 %
第 6 次测试准确率为: 90.0 %
第 7 次测试准确率为: 100.0 %
第 8 次测试准确率为: 86.66666666666667 %
第 9 次测试准确率为: 93.33333333333333 %
第 10 次测试准确率为: 100.0 %
样本数量为: 135
测试数量为: 15
属性类为: 4
标签结果: 3
K值为: 5
测试准确率为: 94.33333333333334 %
```

## 分析:

使用 K 折交叉验证有效避免了过拟合。

## 数据集二: letter-recognition

1	0, 3, 4, 4, 3, 2, 8, 7, 7, 5, 7, 6, 8, 2, 8, 3, 8
2	C, 7, 10, 5, 5, 2, 6, 8, 6, 8, 11, 7, 11, 2, 8, 5, 9
3	T, 6, 11, 6, 8, 5, 6, 11, 5, 6, 11, 9, 4, 3, 12, 2, 4
4	J, 2, 2, 3, 3, 1, 10, 6, 3, 6, 12, 4, 9, 0, 7, 1, 7
5	J, 1, 3, 2, 2, 1, 8, 8, 2, 5, 14, 5, 8, 0, 7, 0, 7
6	H, 4, 5, 5, 4, 4, 7, 7, 6, 6, 7, 6, 8, 3, 8, 3, 8
7	D, 3, 1, 4, 2, 2, 7, 7, 6, 7, 6, 5, 5, 2, 8, 3, 7
8	I, 0, 3, 0, 1, 0, 7, 7, 1, 7, 7, 6, 8, 0, 8, 2, 8
9	H, 2, 3, 4, 2, 2, 8, 7, 3, 5, 10, 6, 8, 3, 8, 2, 8
10	T, 9, 11, 9, 8, 7, 7, 10, 2, 9, 11, 9, 5, 4, 11, 5, 5
11	Y, 7, 9, 7, 7, 4, 3, 10, 2, 7, 11, 11, 7, 1, 11, 2, 5

## 结果:

```
第 1 次测试准确率为: 82.7485380116959 %
第 2 次测试准确率为: 85.08771929824562 %
第 3 次测试准确率为: 87.71929824561403 %
第 4 次测试准确率为: 87.13450292397661 %
第 5 次测试准确率为: 92.10526315789474 %
第 6 次测试准确率为: 89.18128654970761 %
第 7 次测试准确率为: 85.96491228070175 %
第 8 次测试准确率为: 88.30409356725146 %
第 9 次测试准确率为: 86.4406779661017 %
第 10 次测试准确率为: 86.8421052631579 %
样本数量为: 3084
测试数量为: 342
属性类为: 16
标签结果: 26
K值为: 5
测试准确率为: 87.15283972643473 %
```

### 分析:

对于这个数据集, 结果比较理想, 毕竟特征比上一个数据集多不少。

### 数据集三: cmc

1	24, 2, 3, 3, 1, 1, 2, 3, 0, 1
2	45, 1, 3, 10, 1, 1, 3, 4, 0, 1
3	43, 2, 3, 7, 1, 1, 3, 4, 0, 1
4	42, 3, 2, 9, 1, 1, 3, 3, 0, 1
5	36, 3, 3, 8, 1, 1, 3, 2, 0, 1
6	19, 4, 4, 0, 1, 1, 3, 3, 0, 1
7	38, 2, 3, 6, 1, 1, 3, 2, 0, 1
8	21, 3, 3, 1, 1, 0, 3, 2, 0, 1
9	27, 2, 3, 3, 1, 1, 3, 4, 0, 1
10	45, 1, 1, 8, 1, 1, 2, 2, 1, 1

### 结果:



```
第 1 次测试准确率为: 63.94557823129252 %
第 2 次测试准确率为: 31.97278911564626 %
第 3 次测试准确率为: 57.82312925170068 %
第 4 次测试准确率为: 37.41496598639456 %
第 5 次测试准确率为: 55.78231292517006 %
第 6 次测试准确率为: 55.78231292517006 %
第 7 次测试准确率为: 59.183673469387756 %
第 8 次测试准确率为: 63.94557823129252 %
第 9 次测试准确率为: 34.78260869565217 %
第 10 次测试准确率为: 34.26573426573427 %
样本数量为: 1330
测试数量为: 143
属性类为: 9
标签结果: 3
K值为: 4
测试准确率为: 49.48986830974409 %
```

### 分析:

不知道为什么特征比上一个还少, 准确率却低, 想试试 sklearn 里的算法, 无奈总是 import 失败。初步推测是数据的原因吧。

### 回归: 替换核心代码

```
35 def calc(test_x, test_y, inputs, k):
36     # 计算距离
37     temp = [list_sub(temp, inputs) for temp in test_x]
38     temp2 = np.sum(temp, axis=1)
39
40     distance = np.ndarray.tolist(temp2 ** 0.5)
41     # 返回排序后的索引值
42     sorted_dist_indices = arg_sort(distance)
43     class_count = {}
44     # 选择k个最近邻
45     for a in range(k):
46         vote_labels = test_y[sorted_dist_indices[a][0]]
47         temp = str(vote_labels)[1:-1]
48         class_count[temp] = class_count.get(temp, 0) + 1
49     sums = 0
50     for c in class_count:
51         sums = sums + class_count.get(c) * float(c)
52     avg = sums/k
53     return avg
```

```

56 def arg_sort(distance):
57     return sorted(enumerate(distance), key=lambda x: x[1])

```

## 数据集四：Winequality-red

```

1 7.4;0.7;0;1.9;0.076;11;34;0.9978;3.51;0.56;9.4;5
2 7.8;0.88;0;2.6;0.098;25;67;0.9968;3.2;0.68;9.8;5
3 7.9;0.43;0.21;1.6;0.106;10;37;0.9966;3.17;0.91;9.5;5
4 8.5;0.49;0.11;2.3;0.084;9;67;0.9968;3.17;0.53;9.4;5
5 6.9;0.4;0.14;2.4;0.085;21;40;0.9968;3.43;0.63;9.7;6
6 6.3;0.39;0.16;1.4;0.08;11;23;0.9955;3.34;0.56;9.3;5
7 7.6;0.41;0.24;1.8;0.08;4;11;0.9962;3.28;0.59;9.5;5
8 7.9;0.43;0.21;1.6;0.106;10;37;0.9966;3.17;0.91;9.5;5
9 7.1;0.71;0;1.9;0.08;14;35;0.9972;3.47;0.55;9.4;5
10 7.8;0.645;0;2;0.082;8;16;0.9964;3.38;0.59;9.8;6
11 6.7;0.675;0.07;2.4;0.089;17;82;0.9958;3.35;0.54;10.1;5

```

## 结果：

```

99.43714821763602 %
99.62476547842401 %
99.812382739212 %
样本数量为： 1066
测试数量为： 533
属性类为： 11
K值为： 5
总cost为： 316.7999999999996

```

## 分析：

或许是属性有点多，该回归任务完成的不是很理想。对于含有差别较大属性的数据，需要归一化操作。

## 3.1.2、小结

KNN 算法

## 3.2、C4.5 决策树

### 3.2.1、代码实现

```
1  from scipy import *
2  from math import log
3
4
5  # 计算给定数据的香浓熵:
6  def calc_shannon_ent(data_set):
7      num_entries = len(data_set) # 数据总个数
8      label_counts = {} # 标签类别字典 (类别的名称为键, 该类别的个数为值)
9      for featVec in data_set:
10         current_label = featVec[-1] # 选取最后一列, 为标签列
11
12         if current_label not in label_counts.keys(): # 还没添加到字典里的类型
13             label_counts[current_label] = 0
14             label_counts[current_label] += 1
15
16         shannon_ent = 0.0
17         for keys in label_counts: # 求出每种类型的熵
18             prob = float(label_counts[keys]) / num_entries # 每种类型个数占所有的比值
19             shannon_ent -= prob * log(prob, 2)
20
21         return shannon_ent # 返回熵
22
23 # 按照给定的特征划分数据集, 剔除含选出属性的一列, lisan属性为true即为离散
24 def split_data_set(data_set, axis, value, flag, lisan=True):
25     ret_data_set = []
26     if lisan:
27         # 按dataSet矩阵中的第axis列的值等于value的分数据集
28         for featVec in data_set:
29             # 值等于value的, 每一行为新的列表 (去除第axis个数据)
30             if featVec[axis] == value:
31                 reduced_feat_vec = featVec[:axis]
32                 reduced_feat_vec.extend(featVec[axis + 1:])
33                 ret_data_set.append(reduced_feat_vec)
34     else:
35         for featVec in data_set:
```

```

34         if flag:
35             if featVec[axis] >= value:
36                 reduced_feat_vec = featVec[:axis]
37                 reduced_feat_vec.extend(featVec[axis + 1:])
38                 ret_data_set.append(reduced_feat_vec)
39             else:
40                 if featVec[axis] < value:
41                     reduced_feat_vec = featVec[:axis]
42                     reduced_feat_vec.extend(featVec[axis + 1:])
43                     ret_data_set.append(reduced_feat_vec)
44         return ret_data_set # 返回分类后的新矩阵
45
46
47 # 选择最好的数据集划分方式
48 def choose_best_feature_to_split(data_set):
49     num_features = len(data_set[0]) - 1 # 求属性的个数, 减一是因为数据集里有结果那一列
50     base_entropy = calc_shannon_ent(data_set) # 计算 Ent(D) 计算纯度
51     best_infoGain = 0.0 # 最大的信息增益
52     best_feature = -1 # 选择出信息增益最大的列的下标
53     for i in range(num_features): # 求所有属性的信息增益
54         feat_list = [example[i] for example in data_set] # 获得第i列的值

```

```

55     new_entropy = 0.0
56     split_info = 0.0 # 属性的固有值, 属性类别越多, 值越大
57     # 判断是否为离散属性
58     if type(feat_list[0]) is float: # 连续型
59         flag = True # 标记为连续型
60         temp_list = list(set(sorted(feat_list))) # 排序
61         candidate = [] # 候选值
62         # temp_list = list(set(temp_list))
63         for j in range(len(temp_list) - 1):
64             candidate.append((temp_list[j] + temp_list[j + 1]) / 2) # 计算中位点
65         unique_values = set(candidate) # 所有候选值
66         # 考察每一个候选值, 选出最大信息熵所对应的候选值作为阈值
67         infoGain = 0
68         for value in unique_values: # 求第i列属性每个不同值的熵*他们的概率
69             # 获取只含小于value的样本集
70             sub_data_set = split_data_set(data_set, i, value, 0, False)
71             # 获取只含属性大于等于value的样本集
72             sub_data_set1 = split_data_set(data_set, i, value, 1, False)
73             prob_a = len(sub_data_set) / float(len(data_set)) # 求出该值在i列属性中的概率
74             # 计算条件熵
75             new_entropy = \
76                 prob_a * calc_shannon_ent(sub_data_set) + (1 - prob_a) * calc_shannon_ent(sub_data_set1)
77             sub_infoGain = (base_entropy - new_entropy) # 求出第i列属性的信息增益
78             if sub_infoGain > infoGain: # 记录阈值信息

```

```

79         infoGain = sub_infoGain
80         if sub_infoGain > best_infoGain:
81             best_feature = [value, i] # 我需要记录i值
82     else:
83         flag = False
84         unique_values = set(feats_list) # 第i列属性的取值（不同值）数集合
85         for value in unique_values: # 求第i列属性每个不同值的熵*他们的概率
86             # 获取只含属性等于value的样本集
87             sub_data_set = split_data_set(data_set, i, value, 0)
88             # 求出该值在i列属性中的概率
89             prob = len(sub_data_set) / float(len(data_set))
90             # 求i列属性各值对应的熵求和
91             new_entropy += prob * calc_shannon_ent(sub_data_set)
92             split_info -= prob * log(prob, 2)
93             # 求出第i列属性的信息增益率
94         infoGain = (base_entropy - new_entropy) / split_info
95         # 保存信息增益率最大的信息增益率值以及所在的下标（列下标i）
96         if infoGain > best_infoGain:
97             best_infoGain = infoGain
98             if not flag:
99                 best_feature = i
100     return best_feature

```

```

103 # 找出出现次数最多的分类名称
104 def majority_cnt(class_list):
105     class_count = {}
106     for vote in class_list:
107         if vote not in class_count.keys():
108             class_count[vote] = 0
109             class_count[vote] += 1
110     max_1 = 0
111     for d in class_count:
112         if max_1 < class_count.get(d):
113             max_1 = class_count.get(d)
114             result = d
115     return result
116
117 # 创建树
118 def create_tree(data_set, labels_c):
119     # 创建需要创建树的训练数据的结果列表
120     class_list = [example[-1] for example in data_set]
121     # 如果所有的训练数据都是属于一个类别，则返回该类别
122     if class_list.count(class_list[0]) == len(class_list):
123

```

```

124         return class_list[0]
125     # 训练数据只给出类别数据（没给任何属性值数据），返回出现次数最多的分类名称
126     if len(data_set[0]) == 1:
127         return majority_cnt(class_list)
128     if len(data_set) <= 3:
129         return majority_cnt(class_list)
130     # 选择信息增益最大的属性进行分（返回值是属性类型列表的下标，连续型为节点值和下标）
131     best_feat = choose_best_feature_to_split(data_set)
132
133     # 剪枝操作
134     if type(best_feat) == list: # 为连续属性
135         best_feat_label = labels_c[best_feat[1]]
136         my_tree = {best_feat_label: {}} # 以bestFeatLabel为根节点建一个空树
137         del (labels_c[best_feat[1]]) # 从属性列表中删掉已经被选出来当根节点的属性
138         key1 = "<" + str(best_feat[0])
139         key2 = ">=" + str(best_feat[0])
140         # 这里大概花了我两天的时间，因为python list递归会出现问题，
141         # 所以需要复制操作 sub_labels = labels_c[:]
142         sub_labels = labels_c[:]
143         # split_data_set 函数 倒数第二个参数 1代表选中大于value的 0 代表小于
144         return_tree1 = \
145             create_tree(split_data_set(data_set, best_feat[1], best_feat[0], 0, False), sub_labels)
146         sub_labels = labels_c[:]

```

```

147         return_tree2 = \
148             create_tree(split_data_set(data_set, best_feat[1], best_feat[0], 1, False), sub_labels)
149         my_tree[best_feat_label][key1] = return_tree1 # 根据各个分支递归创建树
150         my_tree[best_feat_label][key2] = return_tree2 # 根据各个分支递归创建树
151     else:
152         best_feat_label = labels_c[best_feat] # 根据下表找属性名称当树的根节点
153         my_tree = {best_feat_label: {}} # 以bestFeatLabel为根节点建一个空树
154         del (labels_c[best_feat]) # 从属性列表中删掉已经被选出来当根节点的属性
155         # 找出该属性所有训练数据的值（创建列表）
156         feat_values = [example[best_feat] for example in data_set]
157         unique_values = set(feat_values) # 求出该属性的所有值得集合（集合的元素不能重复）
158         for value in unique_values: # 根据该属性的值求树的各个分支
159             sub_labels = labels_c[:]
160             # 根据各个分支递归创建树
161             my_tree[best_feat_label][value] = \
162                 create_tree(split_data_set(data_set, best_feat, value, 0), sub_labels)
163         return my_tree # 生成的树
164
165     # 实用决策树进行分类
166
167     def classify(input_tree, feat_labels, test_vec):

```

```

169     global class_label
170     first_str = str(input_tree.keys())[12:-3] # 得到第一个字典的key
171     second_dict = input_tree[first_str]
172     feat_index = feat_labels.index(first_str)
173     global class_label
174     for key1 in second_dict.keys():
175         if key1.count("<"): # 连续
176             key2 = float(key1[1:])
177             if test_vec[feat_index] <= key2:
178                 if type(second_dict[key1]).__name__ == 'dict': # 如果没有到最底层, 递归
179                     class_label = classify(second_dict[key1], feat_labels, test_vec)
180                 else:
181                     class_label = second_dict[key1]
182             elif key1.count(">="):
183                 key2 = float(key1[2:])
184                 if test_vec[feat_index] > key2:
185                     if type(second_dict[key1]).__name__ == 'dict': # 如果没有到最底层, 递归
186                         class_label = classify(second_dict[key1], feat_labels, test_vec)
187                     else:
188                         class_label = second_dict[key1]
189             else: # 离散
190                 if test_vec[feat_index] == key1:
191                     if type(second_dict[key1]).__name__ == 'dict': # 如果没有到最底层, 递归
192                         class_label = classify(second_dict[key1], feat_labels, test_vec)
193                 else:
194                     class_label = second_dict[key1]
195     return class_label
196
197 # 读取数据文档中的训练数据 (生成二维列表)
198 def create_train_data():
199     lines_set = open('E:\myCodes\python\mlTest\venv\src\data\letter_train.txt').readlines()
200     data_x = []
201
202     for temp in lines_set:
203         # 将标签移至最后一列
204         temp2 = list(temp[:].split("\n")[0].split(","))
205         # print(temp2)
206         c = temp2[-1]
207         temp2[-1] = temp2[0]
208         temp2[0] = c
209         temp2 = list(map(float, temp2[:-1])) + temp2[-1:]
210         # print("##", temp2)
211         data_x.append(temp2)
212     label = ['A1', 'A2', 'A3', 'A4', 'A5', 'A6', 'A7', 'A8', 'A9', 'A10', 'A11',
213             'A12', 'A13', 'A14', 'A15', 'A16']
214     return data_x, label
215
216 # 读取数据文档中的测试数据 (生成二维列表)
217 def create_test_data():
218     lines_set = open('E:\myCodes\python\mlTest\venv\src\data\letter-test.txt').readlines()
219     data_x = []
220     data_y = []
221     for temp in lines_set:
222         temp2 = list(temp[:].split("\n")[0].split(","))
223         c = temp2[-1]

```

```

226         temp2[-1] = temp2[0]
227         temp2[0] = c
228         temp3 = list(map(float, temp2[:-1]))
229         data_x.append(temp3)
230         data_y.append(temp2[-1:])
231     return data_x, data_y
232
233
234 ▶ if __name__ == '__main__':
235     my_dat_x, labels = create_train_data()
236     myTree = create_tree(my_dat_x, labels)
237     print(myTree)
238     # 重新创建一次属性列表，具体原因尚不太明白
239
240     boot_labels = ['A1', 'A2', 'A3', 'A4', 'A5',
241                   'A6', 'A7', 'A8', 'A9', 'A10', 'A11',
242                   'A12', 'A13', 'A14', 'A15', 'A16']
243     test_list, test_labels = create_test_data()

```

```

244     labels_num = {}
245     total = len(test_labels)
246     pre = []
247     count = 0
248     for testData in test_list:
249         pre.append(classify(myTree, boot_labels, testData))
250     for ic in range(total):
251         key = set(test_labels[ic])
252         labels_num[str(key)] = 1
253         if pre[ic] == str(test_labels[ic])[2:-2]:
254             count += 1
255     print("样本数量为:", len(my_dat_x))
256     print("测试数量为:", len(test_list))
257     print("属性类为:", len(my_dat_x[0])-1)
258     print("标签结果:", len(labels_num))
259     print("测试准确率为:", count / total * 100, "%")

```

C4.5 决策树的代码花了我很多时间，先完成对离散属性的分类，再完成离散与连续属性的结合。其中辨别离散还是连续属性看属性如果是 float 型则为连续，如果是字符串，则是离散，这就要求事先对数据集进行处理到位。



## 数据集一：Car Evaluation

### 离散属性

```
1 vhigh,high,3,more,big,med,unacc
2 vhigh,high,3,more,big,high,unacc
3 vhigh,high,4,2,small,low,unacc
4 vhigh,vhigh,2,2,small,low,unacc
5 vhigh,vhigh,2,2,small,med,unacc
6 vhigh,vhigh,2,2,small,high,unacc
7 vhigh,vhigh,2,2,med,low,unacc
8 vhigh,vhigh,2,2,med,med,unacc
9 vhigh,vhigh,2,2,med,high,unacc
10 vhigh,vhigh,2,2,big,low,unacc
11 vhigh,vhigh,2,2,big,med,unacc
12 vhigh,vhigh,2,2,big,high,unacc
```

### 结果：

```
第 1 次测试准确率为： 81.97674418604652 %
第 2 次测试准确率为： 68.6046511627907 %
第 3 次测试准确率为： 75.5813953488372 %
第 4 次测试准确率为： 95.34883720930233 %
第 5 次测试准确率为： 69.18604651162791 %
第 6 次测试准确率为： 80.23255813953489 %
第 7 次测试准确率为： 81.97674418604652 %
第 8 次测试准确率为： 70.93023255813954 %
第 9 次测试准确率为： 83.13953488372093 %
第 10 次测试准确率为： 64.0625 %
样本数量为： 1600
测试数量为： 128
属性类为： 6
标签结果： 4
测试平均准确率为： 77.10392441860466 %
```

**分析：**使用了简单的预剪枝操作，即如果样本小于等于 3 则不继续进行分类。

## 数据集二：letter-recognition

### 连续属性

```
459 0,3,4,4,3,2,8,7,7,5,7,6,8,2,8,3,8
460 C,7,10,5,5,2,6,8,6,8,11,7,11,2,8,5,9
461 T,6,11,6,8,5,6,11,5,6,11,9,4,3,12,2,4
462 J,2,2,3,3,1,10,6,3,6,12,4,9,0,7,1,7
463 J,1,3,2,2,1,8,8,2,5,14,5,8,0,7,0,7
464 H,4,5,5,4,4,7,7,6,6,7,6,8,3,8,3,8
465 S,3,2,3,3,2,8,8,7,5,7,5,7,2,8,9,8
466 O,6,11,7,8,5,7,6,9,6,7,5,9,4,8,5,5
467 J,3,6,4,4,2,6,6,4,4,14,8,12,1,6,1,6
468 C,6,11,7,8,3,7,8,7,11,4,7,14,1,7,4,8
469 M,7,11,11,8,9,3,8,4,5,10,11,10,10,9,5,7
```

### 结果：

```
第 1 次测试准确率为： 67.2514619883041 %
第 2 次测试准确率为： 71.05263157894737 %
第 3 次测试准确率为： 62.28070175438597 %
第 4 次测试准确率为： 67.54385964912281 %
第 5 次测试准确率为： 66.66666666666666 %
第 6 次测试准确率为： 69.00584795321637 %
第 7 次测试准确率为： 65.2046783625731 %
第 8 次测试准确率为： 66.66666666666666 %
第 9 次测试准确率为： 65.2046783625731 %
第 10 次测试准确率为： 69.2982456140351 %
样本数量为： 3084
测试数量为： 342
属性类为： 16
标签结果： 26
测试平均准确率为： 67.01754385964912 %
```

## 数据集三：winequality-red

### 连续属性

```
7.3;0.65;0;1.2;0.065;15;21;0.9946;3.39;0.47;10;7
7.8;0.58;0.02;2;0.073;9;18;0.9968;3.36;0.57;9.5;7
7.5;0.5;0.36;6.1;0.071;17;102;0.9978;3.35;0.8;10.5;5
6.7;0.58;0.08;1.8;0.097;15;65;0.9959;3.28;0.54;9.2;5
7.5;0.5;0.36;6.1;0.071;17;102;0.9978;3.35;0.8;10.5;5
5.6;0.615;0;1.6;0.089;16;59;0.9943;3.58;0.52;9.9;5
7.8;0.61;0.29;1.6;0.114;9;29;0.9974;3.26;1.56;9.1;5
8.9;0.62;0.18;3.8;0.176;52;145;0.9986;3.16;0.88;9.2;5
8.9;0.62;0.19;3.9;0.17;51;148;0.9986;3.17;0.93;9.2;5
8.5;0.28;0.56;1.8;0.092;35;103;0.9969;3.3;0.75;10.5;7
```

### 结果：

```
第 1 次测试准确率为： 58.0952380952381 %
第 2 次测试准确率为： 51.42857142857142 %
第 3 次测试准确率为： 44.761904761904766 %
第 4 次测试准确率为： 54.285714285714285 %
第 5 次测试准确率为： 43.39622641509434 %
第 6 次测试准确率为： 56.19047619047619 %
第 7 次测试准确率为： 37.142857142857146 %
第 8 次测试准确率为： 48.57142857142857 %
第 9 次测试准确率为： 40.0 %
第 10 次测试准确率为： 33.33333333333333 %
样本数量为： 948
测试数量为： 105
属性类为： 11
标签结果： 5
测试平均准确率为： 46.720575022461816 %
```

**分析：**感觉准确率有点低，但是我使用 python 的 sklearn 库，也是只有 48%说明我的算法问题不大。一时半会也找不出问题所在。

数据集四：Fertility

离散+连续属性：

属性信息：

进行分析的季节。1) 冬天, 2) 春天, 3) 夏天, 4) 秋天。(-1, -0.33, 0.33, 1)

分析时的年龄。18-36 (0,1)

幼稚疾病（即水痘，麻疹，腮腺炎，小儿麻痹症）1) 是, 2) 否。(0,1) 事故或严重创伤 1) 是, 2) 否。(0,1)

手术干预 1) 是, 2) 否。(0,1)

去年高烧 1) 不到三个月前, 2) 三个多月前, 3) 没有。(-1,0,1)

饮酒次数 1) 每天几次, 每天 2 次, 每周 3 次, 每周几次, 每周一次, 每周一次, 几乎永远或从不 (0,1)

吸烟习惯 1) 从不, 2) 偶尔 3) 每天。(-1,0,1)

每天坐着的小时数 ene-16 (0, 1)

输出：诊断正常 (N)，改变 (O)

1	-0.33, 0.69, 0, 1, 1, 0, 0.8, 0, 0.88, N
2	-0.33, 0.94, 1, 0, 1, 0, 0.8, 1, 0.31, 0
3	-0.33, 0.5, 1, 0, 0, 0, 1, -1, 0.5, N
4	-0.33, 0.75, 0, 1, 1, 0, 1, -1, 0.38, N
5	-0.33, 0.67, 1, 1, 0, 0, 0.8, -1, 0.5, 0
6	-0.33, 0.67, 1, 0, 1, 0, 0.8, 0, 0.5, N
7	-0.33, 0.67, 0, 0, 0, 0, -1, 0.8, -1, 0.44, N
8	-0.33, 1, 1, 1, 1, 0, 0.6, -1, 0.38, N
9	1, 0.64, 0, 0, 1, 0, 0.8, -1, 0.25, N
0	1, 0.61, 1, 0, 0, 0, 1, -1, 0.25, N

结果：

```
第 1 次测试准确率为: 80.0 %  
第 2 次测试准确率为: 80.0 %  
第 3 次测试准确率为: 80.0 %  
第 4 次测试准确率为: 90.0 %  
第 5 次测试准确率为: 100.0 %  
第 6 次测试准确率为: 80.0 %  
第 7 次测试准确率为: 80.0 %  
第 8 次测试准确率为: 80.0 %  
第 9 次测试准确率为: 80.0 %  
第 10 次测试准确率为: 90.0 %  
样本数量为: 90  
测试数量为: 10  
属性类为: 9  
标签结果: 1  
测试平均准确率为: 84.0 %
```

### 分析:

当离散属性里面只有一个取值时, 会导致分裂属性为 1, 导致代码出现除 0 bug, 不过已经修复。

### 3.2.2、小结

生成离散属性决策树代码是参考了网上的一些代码, 我觉得写的很不错, 在我完全读懂他的代码后, 自己补全了关于连续属性的生成结点的部分, 路上遇到的问题很多。

① python 在进行递归时应该使用 `deepcopy` 来保存 list 等变量, 不然你的变量在递归中会出现问题。

② 关于对连续属性划分的修改后的信息增益, 一开始理解错了, 导致代码一直达不到理想的结果。但在一次午饭过后, 我随意注释了一段代码, 结果居然对了, 而后我才慢慢发觉是自己理解错了公式。对连续属性进行划分不需要除以  $IV(a)$ , 即属性  $a$  的固有属性。

③ 没怎么实现预剪枝, 因为我发现在递归途中获取不到已经生成的决策树。所

以无法进行预剪枝操作。对后剪枝理解不是很透彻。

### 3.3、聚类算法 k-means

#### 3.3.1、代码实现

```
1 import numpy as np
2 import random
3
4
5 def distant_calc(centre, data_set):
6     return np.sqrt(np.sum(np.power(centre - data_set, 2)))
7
8
9 def create_test_data():
10     lines_set = open('E:\myCodes\python\mlTest\ve'
11                     'nv\src\data\iris.data.txt').readlines()
12     data_x = []
13     data_y = []
14     for data in lines_set:
15         temp = data.split("\n")[0].split(",")
16         data_x.append(list(map(float, temp[:-1])))
17         c = str(temp[-1:])[2:-2]
18         data_y.append(c)
19     return data_x, data_y
20
21
22 # def create_test_data():
23 #     lines_set = open('E:\myCodes\python\mlTest\ve
24 # nv\src\data\Wholesale customers data.csv').readlines()
```

```

25     #     data_x = []
26     #     for data in lines_set:
27     #         temp = data.split("\n")[0].split(",")
28     #         data_x.append(list(map(float, temp)))
29     #     return data_x
30
31
32     # 初始化聚簇中心，选取K个
33     def init_centre(data_set, k):
34         n = np.shape(data_set)[1]
35         centre = np.mat(np.zeros((k, n)))
36         # print(np.shape(data_set)[0])
37         r = random.sample(list(range(np.shape(data_set)[0])), k)
38         print('初始化簇中心', r)
39         # 随机选四个数据
40         for row in range(k):
41             centre[row] = data_set[r[row]]
42     return centre
43
44
45     def k_means(dataSet, k, distMeans = distant_calc, createCent = init_centre):
46         m = np.shape(dataSet)[0] # 获取数据个数
47         # cluster_infor第一列存放该数据所属的中心点，第二列是该数据到中心点的距离
48         cluster_infor = np.mat(np.zeros((m, 2))) # 用于存放该样本属于哪类及质心距离

```

```

49         # print(data_set)
50         centroids = createCent(dataSet, k) # 初始化类簇
51         cluster_changed = True # 用来判断聚类是否已经收敛
52         while cluster_changed:
53             cluster_changed = False
54             for i in range(m): # 把每一个数据点划分到离它最近的中心点
55                 min_dist = float('inf')
56                 min_index = -1
57                 for j in range(k):
58                     distJI = distMeans(centroids[j], np.mat(dataSet[i]))
59                     if distJI < min_dist:
60                         min_dist = distJI
61                         # 如果第i个数据点到第j个中心点更近，则将i归属为j
62                         min_index = j
63                 if cluster_infor[i, 0] != min_index:
64                     cluster_changed = True # 如果分配发生变化，则需要继续迭代
65                     # 并将第i个数据点的分配情况存入字典
66                     cluster_infor[i, :] = min_index, min_dist

```

```

67         for cent in range(k): # 重新计算中心点
68             # 去第一列等于cent的所有列
69             new_clust = dataSet[np.nonzero(cluster_infor[:, 0].A == cent)[0]]
70             # 算出这些数据的中心点
71             centroids[cent, :] = np.mean(new_clust, axis=0)
72     return centroids, cluster_infor
73
74
75     # 使用外部指标度量聚类性能
76 def outer_judge(result, models):
77     # 获取|SS|
78     count_a = 0
79     count_b = 0
80     count_c = 0
81     count_d = 0
82     n = len(models)
83     for i in range(n):
84         for j in range(i+1, n):
85             if result[i].tolist()[0][0] \
86                 == result[j].tolist()[0][0] and models[i] == models[j]:
87                 count_a += 1
88             elif result[i].tolist()[0][0] \
89                 == result[j].tolist()[0][0] and models[i] != models[j]:
90                 count_b += 1
91
92             elif result[i].tolist()[0][0] \
93                 != result[j].tolist()[0][0] and models[i] == models[j]:
94                 count_c += 1
95             else:
96                 count_d += 1
97
98     JC = count_a / (count_a + count_b + count_c)
99     FMI = np.sqrt((count_a / (count_a +
100                         count_b)) * (count_a / (count_a + count_c)))
101     Rand = 2 * (count_a + count_d) / (n * (n-1))
102     print('JC系数为: ', JC)
103     print('FM指数为: ', FMI)
104     print('Rand指数为: ', Rand)
105
106 def inner_judge(result, data_set, centroids):
107     n = np.shape(result)[0]
108
109     # 用于统计每个类簇中样本所对应应在原始数据里的下标
110     dic_cluster = {}

```



```

110     for i in range(n):
111         list1 = []
112         item = result[i].tolist()[0][0]
113         if item not in dic_cluster.keys():
114             dic_cluster[item] = list1
115         else:
116             list1 = dic_cluster[item]
117             list1.append(i)
118     dic_avg = {}
119     for i in dic_cluster:
120         total_diatant = 0
121         tmp = dic_cluster[i]
122         for j in range(len(tmp)):
123             for k in range(j + 1, len(tmp)):
124                 total_diatant += \
125                     distant_calc(np.mat(data_set[tmp[j]]), np.mat(data_set[tmp[k]]))
126         dic_avg[i] = (total_diatant * 2 / (n * (n - 1)))
127
128     # 计算各个中心点距离
129     num_cluster = len(dic_avg)
130     total_max = 0
131     for i in range(num_cluster):
132         max_distant = 0
133         for j in range(num_cluster):
134             new_distant = (dic_avg[i] + dic_avg[j]) \
135                 / distant_calc(np.mat(centroids[i]), np.mat(centroids[j]))
136             if new_distant > max_distant:
137                 max_distant = new_distant
138             total_max += max_distant
139     DBI = total_max / num_cluster
140     print('DB指数为: ', DBI)
141     return DBI
142
143
144
145 if __name__ == '__main__':
146     k = 3
147     # data_set = create_test_data()
148     data_set, labels = create_test_data()
149     centroids, result = k_means(np.mat(data_set), k)
150     outer_judge(result, labels)
151     # min_DBI = 1
152     # for i in range(3, 15):
153     #     centroids, result = k_means(np.mat(data_set), i)
154     #     DBI = inner_judge(result, data_set, centroids)
155     #     if DBI < min_DBI:
156         min_DBI = DBI

```

```

157     #         k = i
158
159     print('K值为:', k)
160     print('簇坐标', centroids)
161     print('数据总数', len(data_set), '已分类数据:', np.shape(result)[0])
162     dic = {}
163     for i in result:
164         item = i.tolist()[0][0]
165         if item not in dic.keys():
166             dic[item] = 0
167         dic[item] += 1
168     print('分类结果', dic)
169
170

```

## 数据集一：Iris

使用外部指标度量

**结果：**

```

↓      初始化簇中心 [116, 47, 92]
↕      JC系数为： 0.6998950682056663
↕      FM指数为： 0.8235798153491702
↕      Rand指数为： 0.8818181818181818
↕      簇坐标 [[6.81666667 3.06333333 5.62666667 2.03          ]
[4.98333333 3.38333333 1.46428571 0.24761905]
[5.89591837 2.75510204 4.39183673 1.44081633]]
↕      数据总数 121 已分类数据： 121
↕      分类结果 {1.0: 42, 2.0: 49, 0.0: 30}

```

**分析：**

聚类性能度量外部指标 JC 系数、FM 指数、Rand 指数取值在[0,1],均是越大越好，说明聚类效果还是不错的。另外，在进行多次测试后发现，有一定的几率出现很大的偏差。我猜测这和初始化簇中心有关。

数据集二：winequality-red

**结果：**

```
初始化簇中心 [143, 541, 860, 817, 118, 476]
JC系数为： 0.1618039108222905
FM指数为： 0.29612237485791715
Rand指数为： 0.6134293111479423
K值为： 6
特征数： 6
数据总数 1053 已分类数据： 1053
分类结果 {1.0: 189, 0.0: 208, 4.0: 244, 3.0: 106, 2.0: 137, 5.0: 169}
```

**分析：**

效果没有上一个数据集好，特别是 JC 系数太小。同一个数据集，之前使用 c4.5 决策树的时候也发现这个数据集有点问题，现在可以下定论，这个数据集确实存在一些问题。

数据集三：seeds

使用内部指标度量

**结果：**

**K 取值为 3：**

```
初始化簇中心 [108, 33, 44]
DB指数为： 0.022658957675026903
簇坐标 [[18.72180328 16.29737705 0.88508689 6.20893443 3.72267213 3.60
6.06609836 1.98360656]
[11.90906667 13.25026667 0.85154933 5.22233333 2.86509333 4.72218667
5.09304 2.86666667]
[14.63202703 14.45324324 0.8790973 5.56178378 3.27489189 2.74404324
5.18493243 1.13513514]]
数据总数 210 已分类数据： 210
分类结果 {2.0: 74, 1.0: 75, 0.0: 61}
```

## K 取值为 12:

```
↓ 初始化簇中心 [116, 96, 135, 53, 28, 148, 125, 78, 61, 95, 46, 173]
DB指数为: 0.0003163842035344358
簇坐标 [[19.15875    16.495      0.8848625    6.295625    3.7756875
        6.152375    2.          ]
 [18.80111111    16.30055556    0.88917222    6.19722222    3.72533333    2.24
        6.08405556    2.          ]
 [15.56833333    14.91416667    0.879775     5.72325     3.38408333    3.82
        5.48075     1.5          ]
```

## 分析:

采用了聚类性能度量内部指标 DBI 来度量结果, 可见, 取不同的 k 对结果的影响相当大, 一般的 k 值在 3 到 10, 可以采用逐个尝试的方式来确定 k 值。随着 K 值的增加, 内部指标 DBI 肯定是越来越小, 这里我们需要设定一个阈值, 直到 DBI 小于阈值, 我们就认为该聚类算法已经完成任务。刚开始对 DBI 的理解出了点偏差, 试了很多 k 值, 都是取最大 k 值时 DBI 最小, 当时有点想不通, 以为是代码的问题, 闹了个笑话。

## 3.3.2、小结

- ① 在初始化簇中心生成随机数时, 一开始使用的 `random.randint()` 函数出现了重复数据, 后面改用 `random.sample()` 函数解决了这个问题。
- ② 使用内部度量指标应该设定阈值。

## 3.4、朴素贝叶斯分类器

### 3.4.1、代码实现

```
1  import copy
2  from math import log
3  import random
4
5
6  # 定义一个标签类，用于存储标签信息
7  class Lable(object):
8      def setValue(self, value):
9          self.value = value
10
11      def setIdx_list(self, idx_list):
12          self.idx_list = idx_list
13
14      def setCount(self, count):
15          self.count = count
16
17
18  # 读取数据文档中的训练数据（生成二维列表）
19  def load_train_data(k):
20      lines_set = open('E:\myCodes\python\mlTest\venv\src\d'
21                      'ata\PhishingData.arff.txt').readlines()
22      data = []
23      for temp in lines_set:
24          temp2 = list(temp[:].split("\n")[0].split(","))
25          data.append(temp2[:])
26      list1 = list(range(0, len(data), k))
27      w = random.sample(list1, 1)[0]
28      data_x = []
29
30      data_x.extend(data[:w])
31      data_x.extend(data[w + int(len(data) / k):])
32      data_y = [str(i[-1:])[1:-1] for i in data_x]
33      data_x = [i[:-1] for i in data_x]
34      test_y = [str(i[-1:])[1:-1] for i in data[w:w +
35                                                    int(len(data) / k)]]
36      test_x = [i[:-1] for i in data[w:w + int(len(data) / k)]]
37      return data_x, data_y, test_x, test_y
```

```

49
50
51 def init_table(data_set, lables):
52     n = len(data_set)
53     lab_list = []
54     lab_set = set(lables)
55     # 向标签类存数据
56     for lab in lab_set:
57         new_lab = Lable()
58         new_lab.setValue(lab)
59         count = 0
60         sub_lab_list = []
61         for index_lab in range(n):
62             if lables[index_lab] == lab:
63                 count += 1
64                 sub_lab_list.append(index_lab)
65         new_lab.setCount(count)
66         new_lab.setIdx_list(sub_lab_list)
67         lab_list.append(new_lab)
68     list_resu = []
69     # 遍历每个属性
70     for i in range(len(data_set[0])):

```

```

71         attra_i = [] # 存放所有样本在第i个属性的所有值
72         # 遍历某个属性的所有样本
73         for j in range(n):
74             attra_i.append(data_set[j][i]) # j行, i列
75         unique_att = set(attra_i)
76         dic_att = {}
77
78         for att in unique_att:
79             dic_value = {} # 存放属性值att所对应所有标签的样本个数
80
81             # 遍历每个标签
82             for item in lab_list:
83                 lable = item.value # value即是标签
84                 for inner in item.idx_list:
85                     if data_set[inner][i] == att:
86                         if lable not in dic_value.keys():
87                             dic_value[lable] = 0
88                             dic_value[lable] += 1
89

```

```

90         # 根据计数结果计算概率
91         if label in dic_value.keys():
92             dic_value[label] = \
93                 (dic_value[label] + 1) / (item.count + len(unique_att))
94         # 将未出现的属性处理
95         for la in lab_set:
96
97             if la not in dic_value.keys():
98                 dic_value[la] = 1 / len(data_set)
99             dic_att[att] = copy.deepcopy(dic_value)
100             list_resu.append(copy.deepcopy(dic_att))
101         return list_resu
102
103
104     def native_bayes(test_data, table_prob):
105         n = len(test_data)
106         list2 = [] # 存放结果
107         list_tmp = []
108         list_prob = [] # 存放每一个样本对应每个标签的概率的集合
109         for item in range(n):
110             # 第i个属性， 逐个计算每一个属性对应的概率
111             for i in range(len(test_data[0])):
112                 value = test_data[item][i] # 代表第item行，i列的值
113                 list_prob.append(table_prob[i][value])
114             dic_line = {}
115             for sub_pro in list_prob:
116                 for ree in sub_pro:
117
118                     if ree not in dic_line.keys():
119                         dic_line[ree] = 1
120                     tmp = log(sub_pro[ree], 2)
121                     dic_line[ree] += tmp
122                 list_tmp.append(dic_line)
123             for item in list_tmp:
124                 max_prob = float('-inf')
125                 max_key = ''
126                 for dic_item in item:
127                     if max_prob < item[dic_item]:
128                         max_prob = item[dic_item]
129                         max_key = dic_item
130                 list2.append(max_key)
131             return list2
132
133     if __name__ == '__main__':
134         K = 10
135         avg_accuracy = 0
136         for i in range(K):
137             count = 0
138             data_set, lables, test_set, test_lables = \
139                 load_train_data(K)
140             labels_num = set(lables)
141             total = len(test_lables)

```

```

135     table_prob = init_table(data_set, lables)
136     predict = native_bayes(test_set, table_prob)
137     for j in range(len(test_set)):
138         if predict[j] == test_labels[j]:
139             count += 1
140     accuracy = count / total * 100
141     print("第", (i + 1), "次测试准确率为: ", accuracy, "%")
142     avg_accuracy += accuracy / K
143     if i == K-1:
144         print("样本数量为: ", len(data_set))
145         print("测试数量为: ", len(test_set))
146         print("属性类为: ", len(data_set[0]))
147         print("标签结果: ", len(labels_num))
148         print("测试平均准确率为: ", avg_accuracy, "%")

```

## 数据集一：Phishing

### 结果：

```

第 1 次测试准确率为: 62.22222222222222 %
第 2 次测试准确率为: 61.48148148148148 %
第 3 次测试准确率为: 52.59259259259259 %
第 4 次测试准确率为: 56.52173913043478 %
第 5 次测试准确率为: 57.77777777777777 %
第 6 次测试准确率为: 53.48837209302325 %
第 7 次测试准确率为: 60.74074074074074 %
第 8 次测试准确率为: 52.59259259259259 %
第 9 次测试准确率为: 61.48148148148148 %
第 10 次测试准确率为: 48.148148148148145 %
样本数量为: 1218
测试数量为: 135
属性类为: 9
标签结果: 3
测试平均准确率为: 56.70471482604951 %

```

### 使用 sklearn:



```
第 1 次测试准确率为： 85.18518518518519 %  
第 2 次测试准确率为： 77.77777777777779 %  
第 3 次测试准确率为： 85.92592592592592 %  
第 4 次测试准确率为： 78.51851851851852 %  
第 5 次测试准确率为： 81.41592920353983 %  
第 6 次测试准确率为： 78.51851851851852 %  
第 7 次测试准确率为： 82.22222222222221 %  
第 8 次测试准确率为： 81.48148148148148 %  
第 9 次测试准确率为： 78.51851851851852 %  
第 10 次测试准确率为： 85.18518518518519 %  
样本数量为： 1218  
测试数量为： 135  
属性类为： 9  
测试平均准确率为： 81.47492625368731 %
```

### 分析：

说明自己的模型不行，还需要改进，并且我没有使用极大似然法，模型太过简单，只是采用了先验概率来计算。

### 3.4.2、小结

关于朴素贝叶斯分类器，准确率有点低，在检查完代码后，觉得是朴素贝叶斯分类器的条件比较苛刻，sklearn 也只要 80%的准确率，要求各个特征独立，而且相等，还有就是和自己假设的模型也有关，所以我就值测试这一组值了。

## 4、总结

结合这些天对机器学习学习和理解，发现不仅仅是掌握算法那么简单，还需要学会优化其速度和准确度。自己的代码能力还有限，需要借助网上大佬们的代码。机器学习算法有很多，但最重要的是学会每类算法的核心思想。

另外，训练数据的处理和挑选也是一门学问，数据应该选取适合自己算法的。

模型的参数选择也很重要，比如 K-means 算法，直接影响到最终的结果。经验也很重要。

最后，实验过程中发现了一个有问题的数据集，也算是个小乐趣。有时候并不是全是自己的错。