


```

#         LabelCounts[currentLabel] = 0                                # 这一步其实就是
#         LabelCounts[currentLabel] += 1                                # 统计有多少个类以
#     shannonEnt = 0
#     for key in LabelCounts:
#         prob=float(LabelCounts[key]) / numEntries                    # 计算单个类的熵
#         shannonEnt -= prob*log(prob,2)                                # 累加每个类的熵值
#     return shannonEnt

```

In [3]:

```

# #创建数据集
# def createdataSet1():
#     dataSet = [['长', '粗', '男'],
#                 ['短', '粗', '男'],
#                 ['长', '细', '女'],
#                 ['短', '细', '女'],
#                 ['短', '粗', '女'],
#                 ['长', '粗', '女'],]
#     labels = ['头发', '声音'] #特征
#     return dataSet, labels

#*****

```

```

def createDataSet1():    # 创造示例数据
    dataSet = [['长', '粗', '男'],
                ['短', '粗', '男'],
                ['短', '粗', '男'],
                ['长', '细', '女'],
                ['短', '细', '女'],
                ['短', '粗', '女'],
                ['长', '粗', '女'],
                ['长', '粗', '女']]
    features = ['头发', '声音'] #两个特征
    return dataSet, features

```

In [4]:

```

# #分割数据集
# def splitDataSet(dataSet,axis,value):
#     retDataSet = []
#     for featVec in dataSet:
#         if featVec[axis]==value:
#             reducefeatVec = featVec[:axis]
#             reducefeatVec.extend(featVec[axis+1:])
#             retDataSet.append(reducefeatVec)
#     return retDataSet

#*****

def splitDataSet(dataSet,axis,value): # 按某个特征分类后的数据
    retDataSet = []
    for featVec in dataSet:
        if featVec[axis] == value:
            reducedFeatVec = featVec[:axis]
            reducedFeatVec.extend(featVec[axis+1:]) # 这两步的操作是没有包括划分
            retDataSet.append(reducedFeatVec)

    return retDataSet

```

```

In [5]: # #选择最优分类特征
# def chooseBestFeatureToSplit(dataSet):
#     numFeatures = len(dataSet[0])-1
#     baseEntropy = calcSannonEnt(dataSet)
#     bestInfoGain = 0
#     bestFeatture = -1
#     for i in range(numFeatures):
#         featList = [example[i] for example in dataSet]
#         uniqueVals = set(featList)
#         newEntropy = 0
#         for vaule in uniqueVals:
#             subDataSet = splitdataSet(dataSet,i,vaule)
#             prob = len(subDataSet) / float(len(dataSet))
#             newEntropy+=prob * calcSannonEnt(subDataSet)
#             infoGain = baseEntropy-newEntropy
#             if (bestInfoGain<bestInfoGain):
#                 bestFeatture = infoGain
#                 bestFeatture = i
#     return bestFeatture

#*****

def chooseBestFeatureToSplit(dataSet): # 选择最优的分类特征
numFeatures = len(dataSet[0]) - 1 # 获得特征的个数 2个
baseEntropy =calcSannonEnt(dataSet) # 原始的信息熵
bestInfoGain = 0
bestFeature = -1
for i in range(numFeatures): # 遍历两个特征
    featList = [example[i] for example in dataSet]
    uniqueVals = set(featList) # 引入集合
    newEntropy = 0

    for value in uniqueVals:
        subDataSet = splitDataSet(dataSet, i, value) # 根据某个特征分类后的数据集
        prob = len(subDataSet) / float(len(dataSet))
        newEntropy += prob*calcSannonEnt(subDataSet) # 按特征分类后的条件熵
    infoGain = baseEntropy - newEntropy # 原始熵与按特征分类后的熵的差值 即信息增益
    if (infoGain > bestInfoGain): # 若按某特征划分后，熵值减少的最大，则选择该特征
        bestInfoGain = infoGain
        bestFeature = i
return bestFeature # 返回的是最优特征的索引

```

```

In [6]: # def majorityCnt(classList):
#     classCount = {}
#     for vote in classList:
#         if vote not in classList:
#             if vote in classCount.keys():
#                 classCount[vote]=0
#                 classCount[vote]+=1
#             sortedclassCount = sorted(classCount.items(),key=operator.itemgetter(1))
#             print(sortedclassCount)
#     return sortedclassCount[0][0]

#*****

def majorityCnt(classList): # 按分类后类别数量排序，比如：最后分类为2男1女，则返回男
classCount = {}
for vote in classList:
    if vote not in classCount.keys():

```

```

        classCount[vote] = 0
        classCount[vote] += 1
        sortedClassCount = sorted(classCount.items(), key = operator.itemgetter(1),r
        #print(sortedClassCount)
        return sortedClassCount[0][0]

```

```

In [7]: ## 错误示范
def createTree(dataSet,labels):
    classList = [example[-1] for example in dataSet]
    if classList.count(classList[0]) == len(classList):
        return classList[0]
    if len(dataSet[0]) == 1:
        return majorityCnt(classList)
    bestFeat = chooseBestFeatureToSplit(dataSet)
    bestFeatlabel = labels[bestFeat]

    myTree = {bestFeatlabel:{}}
    del(labels[bestFeat])
    featValues = [ example[bestFeat] for example in dataSet]
    uniqueValues = set(featValues)
    for value in uniqueValues:
        subLabels = labels[:]
        myTree[bestFeatlabel][value] = createTree(splitDataSet(dataSet,bestFeat
    return myTree

#####

# 构建决策树(ID3决策树)

# def createTree(dataSet, labels):
#     classList = [example[-1] for example in dataSet]                # 类别:
#     if classList.count(classList[0]) == len(classList):              # 最终叶
#         return classList[0]
#     if len(dataSet[0]) == 1:
#         return majorityCnt(classList)
#     bestFeat = chooseBestFeatureToSplit(dataSet)                    # 选择最
#     bestFeatLabel = labels[bestFeat]
#     myTree = {bestFeatLabel:{}}                                       # 分类结
#     del(labels[bestFeat])                                             # Label
#     featValues = [example[bestFeat] for example in dataSet]
#     uniqueVals = set(featValues)                                       # {'粗
#     for value in uniqueVals:
#         subLabels = labels[:]
#         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFe
#     return myTree

```

```

In [8]: if __name__ == '__main__':
        dataSet, labels = createDataSet1() # 创造示例数据
        print(createTree(dataSet, labels)) # 输出决策树模型结果

```

```
{'声音': {'细': '女', '粗': {'头发': {'短': '男', '长': '女'}}}}
```