## 0.1   Basic Method

# 1   K-Nearest Neightbour Algorithm

The K-Nearest Neightbour Algorithm is an intuitive algorithm. Given an unknown sample, we compute its 'distance'

$$d(\boldsymbol{v_i}, \boldsymbol{v_j}) = \left( \sum_{k=1}^{n} (v_i^{(k)} - v_j^{(k)})^p \right)^{\frac{1}{p}} \tag{1}$$

to elements in our training set, and use the average features of the K nearest elements picked out to predict the unknown one.

## 1.1   Implementation

---
**Algorithm 1** K-NN Algorithm ($S$, $Vec$)

---
**Require:** training set: $S$, sample vector $Vec = \{v_1, v_2, \ldots, v_n\}$, hashmap $H$

   **for** $s_i \in S$ **do**

      $d_i \leftarrow \sqrt{\sum_{j=1}^{n}((s_i)_j - v_j)^2}$

      $H[s_i] = d_i$

   **end for**

   **sort** $H$ with key

   $array \leftarrow s_i | H[s_i]$ are K nearest

   $Vec.label = \sum_{s_i \in array} \epsilon_i \cdot s_i.label, \quad \epsilon_i = H[s_i] / \sum_{s_j \in array} H[s_j]$

---

## 1.2   Performance Analysis

Assume that $\boldsymbol{x}$ is our testing vector, and the nearest training vector is $\boldsymbol{z}$. Then the generalizetion error is

$$P_{\text{err}} = 1 - \sum_{c \in \mathcal{Y}} P(c|\boldsymbol{x})P(c|\boldsymbol{z}) \tag{2}$$

where $\mathcal{Y}$ is the label set. Then assume that the training set is dense enough, such that $\forall\, \boldsymbol{x},\ \exists\, \delta,\ \boldsymbol{z} \in \boldsymbol{x} + \delta$. The condition gives

$$
\begin{aligned}
P_{\text{err}} &= 1 - \sum_{c \in \mathcal{Y}} P(c|\boldsymbol{x})P(c|\boldsymbol{z}) \\
&\approx 1 - \sum_{c \in \mathcal{Y}} P^2(c|\boldsymbol{x}) \\
&= 1 - \left( \arg\max_{c \in \mathcal{Y}} P(c|x) \right)^2 \\
&= \left( 1 - \arg\max_{c \in \mathcal{Y}} P(c|x) \right)\left( 1 + \arg\max_{c \in \mathcal{Y}} P(c|x) \right) \leq 2 - 2\arg\max_{c \in \mathcal{Y}} P(c|x)
\end{aligned}
\tag{3}
$$

This manifests that the error rate of K-NN will not exceed the double of the one of the Bayes optimal classifier.

## 1.3 Code (with Python)

Listing 1: **K-NN.py**

```python
import numpy as np
from numpy import *

def fileToMatrix(filename):
     file =  open(filename)
    arrayOfLines =  file.readlines()
    numOfLines =  len(arrayOfLines)
    returnMat = np.zeros([numOfLines, 3], dtype = double)
    labelVector = []
    index = 0
    for line in arrayOfLines:
        line = line.strip()
        lineList = line.split('\t')
        returnMat[index,...] = lineList[0:3]
        labelVector.append( int(lineList[-1]))
        index += 1
    return returnMat, labelVector

def normalize(dataMat): #normalize the dataset
    colMinVal = dataMat. min(0)
    colMaxVal = dataMat. max(0)
    interval = colMaxVal - colMinVal
    normDataMat = np.zeros(shape(dataMat))
    colLength = shape(dataMat)[0]
    normDataMat = dataMat - np.tile(colMinVal, (colLength, 1))
    np.seterr(invalid = 'ignore')
    normDataMat = np.divide(normDataMat, np.tile(interval, (colLength, 1)))
    return normDataMat

def classify(sampleVec, dataSet, labelVec, K):
    dataSetSize = dataSet.shape[0]
    diffMat = np.tile(sampleVec, (dataSetSize, 1)) - dataSet
    sqDiffMat = diffMat ** 2
    sqDistances = sqDiffMat. sum(axis = 1) #calculate the distance to each element in the dataset
    distances = sqDistances ** (1/2)
    disIndices = distances.argsort() #replace the distances with their ranks
    totalDistance = 0
    dis, vote = {}, {}
```

```
39        for i in  range(dataSetSize):
40            if disIndices[i] < K:
41                dis[disIndices[i]] = (distances[i], i)
42        for d in dis:
43            totalDistance += dis[d][0]
44        for i in dis:
45            weight = dis[i][0] / totalDistance
46            if labelVec[dis[i][-1]] in vote:
47                vote[labelVec[dis[i][-1]]] += weight
48            else:
49                vote[labelVec[dis[i][-1]]] = weight
50        sortedVoteCount =  sorted({v : k for k, v in vote.items()}.items(), reverse = True)
51        return sortedVoteCount[0][1]
52
53    def K_NN(sampleVec, K, dataSet, labelVec):
54        return classify(normalize(sampleVec), normalize(dataSet), labelVec, K)
```

## 1.4   Application: Handwriting Recognition

We use $32 \times 32$ martices to represent the handwritten image, where 1 stands for occupied pixel, and 0 stands for blanket. We have the dataset's dimension deduced to a $1 \times 1024$ vector and be compared with the trained data.



Listing 2: **HWR.py**

```
1   from os import listdir
2   from K_NN import *
3
4   def imageToVec(filename):
5        file =  open(filename)
6       returnVec = np.zeros([1, 1024], dtype =  int)
```

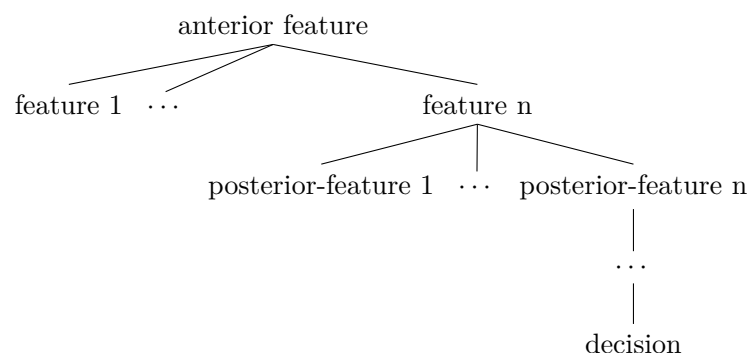```
7        for i in  range(32):
8            lineStr =  file.readline()
9            for j in  range(32):
10               returnVec[0, i * 32 + j] =  int(lineStr[j])
11       return returnVec
12
13   def handWritingRecognition(filename, dataDir): #The training dataset is stored in ./trainingDigits
14       sampleVec = imageToVec(filename)
15       trainingFileList = listdir(dataDir)
16       listLength =  len(trainingFileList)
17       trainingMat = np.zeros([listLength, 1024])
18       labelVec = []
19       for i in  range(listLength):
20           fileName = trainingFileList[i]
21           fileStr = fileName.split('.')[0]
22           fileClass =  int(fileStr.split('_')[0])
23           fileVec = imageToVec('trainingDigits/{}'. format(fileName))
24           trainingMat[i,...] = fileVec
25           labelVec.append(fileClass)
26       return K_NN(sampleVec,  int( input()), trainingMat, labelVec)
```

# 2  Decision Tree

## 2.1  Basic Model

We tend to enable machines to do decision-making like humans. The data structure we use is the decision tree, where each internal node corresponds to a characteristic testing $a_i$, and each leaf node denotes a final decision $y_i$. The core manipulation is to build up optimal classification at each node. Every top-down process on analyzing a sample corresponds to a testing sequence.

## 2.2   Partition Scenario

**Shannon Entropy**

For a decision set $D$, the **information size** $H_0(D)$ which denotes the number of bits needed to encode elements in $D$ is $H_0(D) = \log_2|D|$. Let $\boldsymbol{D} = (D, p)$ be a discrete probability space, where $D = \{D_1, D_2, \ldots, D_n\}$ is a finite set, with $D_i$ corresponds to probability $p_i$ under definite discrete characteristic. Then the **Shannon entropy** of $\boldsymbol{D}$ is
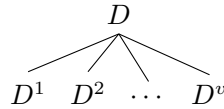
$$\text{Ent}(D) = -\sum_{i=1}^{n} p_i \log_2 p_i \tag{4}$$

Since $-\log_2 x$ is convex, we give an upper-bound for the entropy where

$$-\log_2(\sum_{i=1}^{n} \frac{1}{p_i} p_i) \leq \sum_{i=1}^{n} p_i(-\log_2 \frac{1}{p_i}) = -\text{Ent}(D) \longrightarrow \text{Ent}(D) \leq \log_2 n \tag{5}$$

**Information Gain**

Given a partition criterion $a = \{a^1, a^2, \ldots, a^v\}$ where $a^i \in a$ is the possible value. For sample set $D$, we split it into subsets $D^1, D^2, \ldots, D^v$ where $D^i$ is the subset determined by criterion $a^i$.

$$
\begin{array}{c}
D \\
\diagup \ | \ | \ \diagdown \\
D^1 \quad D^2 \quad \cdots \quad D^v
\end{array}
$$

Then the **information gain** we obatin from this partition is

$$\text{Gain}(D, a) = \text{Ent}(D) - \sum_{i=1}^{v} \frac{|D^i|}{|D|} \text{Ent}(D^i) \tag{6}$$

In ID3 Algorithm, the **optimal class partition** $a_*$ for $A$ with sample $D$ is defined as

$$a_* = \arg\max_{a \in A} \text{Gain}(D, a) \tag{7}$$

while for C4.5 Algorithm, the optimal one is defined as

$$a_* = \arg\max_{a \in A} \text{GainRatio}(D, a), \quad \text{GainRatio}(D, a) = \text{Gain}(D, a)/\text{Ent}(D) \tag{8}$$

# 3   Implementation

Assume that training sample set $D = \{(\boldsymbol{x_1}, y_1), (\boldsymbol{x_2}, y_2), \ldots, (\boldsymbol{x_n}, y_n)\}$, where $\boldsymbol{x_i}$ is the class characteristics, and $y_i$ is the decision, eventually treated as leafnode. Possible partition criterion for $D$ is $A = \{a_1, a_2, \ldots, a_d\}$, where $a_i$ denotes a possible partition.

---
**Algorithm 2** treeGenerate($D$, $A$)
---
**Require:** Training set $D$, Partition criterion set $A$

**initialize** $node$

**if** $\forall D_i D_j \in D, i \neq j, \ D_i = D_j$ **then**

    $node = leafnode, \ node \leftarrow D.y$

    **return**

**end if**

**if** $A = \emptyset$ **or** $\forall a_i, a_j \in A, i \neq j, \ \text{Gain}(D, a_i) = \text{Gain}(D, a_j)$ **then**

    $node = leafnode, \ node \leftarrow \arg\max_{y}(|N|, \ N = \{y | (\boldsymbol{x}, y) \in D\})$

**end if**

$a_* \leftarrow \arg\max_{a \in A} \text{Gain}(D, a)$

**for** $a_*^v \in a_*$ **do**

    **initialize** $node.branch^v$, $D_v$ be the subset splited with $a_*^v$

    **if** $D_v = \emptyset$ **then**

        $node.branch^v = leafnode, node.branch^v \leftarrow \arg\max_{y}(|N|, \ N = \{y | (\boldsymbol{x}, y) \in D_v\})$

    **else**

        $node.branch^v = \text{treeGenerate}(D_v, \ A - \{a_*\})$

    **end if**

**end for**

# 4 Code

Listing 3: **decisionTree.py**

```python
from math import log

class decisionNode( object): # tree organized by decision tree data structure
    def __init__(self, label):
        self.label = label
        self.branches = []
        self.decision = ""
    def assignDecision(self, decision):
        self.decision += decision
    def addBranch(self, newNode):
        newNode.assignDecision(self.decision)
        self.branches.append(newNode)
    def visualize(self, treeNode, layer): # Visualize manipulation displays the layer an label
         belongs to and its anterior choice
        print("({}){}: {}". format(layer, treeNode.decision, treeNode.label))
        for s in treeNode.branches:
            self.visualize(s, layer + 1)

def shannonEntropy(dataset):
    entriesNum =  len(dataset)
```

```python
20        labelCount = {}
21        for dataVec in dataset:
22            dataLabel = dataVec[-1] # The last element in the vector is our decision
23            if dataLabel not in labelCount:
24                labelCount[dataLabel] = 0
25            labelCount[dataLabel] += 1
26        shannonEntropy = 0.000
27        for label in labelCount:
28            probability = labelCount[label] / entriesNum
29            shannonEntropy -= probability * log(probability, 2)
30        return shannonEntropy
31
32    def splitDataset(dataset, index, expectValue): # Search the dataset with specified index and return
            the reduced dataset
33        retDataset = []
34        for lineVec in dataset:
35            if lineVec[index] == expectValue:
36                reducedVec = lineVec[:index]
37                reducedVec += lineVec[index+1:]
38                retDataset.append(reducedVec)
39        return retDataset
40
41    def optimalPartition(dataset): # ID3 Algotihm
42        featureNums =  len(dataset[0]) - 1
43        originalEntropy = shannonEntropy(dataset)
44        bestFeature = -1
45        maxInfoGain = 0.00
46        for featureIndex in  range(featureNums):
47            labelVec =  set([dataset[i][featureIndex] for i in  range( len(dataset))])
48            extraEntropy = 0.00
49            for label in labelVec:
50                reducedSet = splitDataset(dataset, featureIndex, label)
51                extraEntropy +=  len(reducedSet) /  float( len(dataset)) * shannonEntropy(reducedSet)
52            if originalEntropy - extraEntropy > maxInfoGain:
53                maxInfoGain = originalEntropy - extraEntropy
54                bestFeature = featureIndex
55        return bestFeature
56
57    def majorityCount(classList):
58        classNum = {}
59        for member in classList:
60            if member not in classNum:
61                classNum[member] = 0
62            classNum[member] += 1
63        reverseDict = {v:k for k, v in classNum.items()}
```

```
64      orderList =  sorted(reverseDict)
65      return reverseDict[ max(orderList)]
66
67
68  def createTree(dataset, labels):
69      classList = [data[-1] for data in dataset]
70      if classList.count(classList[0]) ==  len(classList):
71          return decisionNode(classList[0])
72      if  len(dataset[0]) == 1:
73          return decisionNode(majorityCount(classList))
74      bestPartition = optimalPartition(dataset)
75      bestPartitionLabel = labels[bestPartition]
76      newTree = decisionNode(bestPartitionLabel)
77      uniqueVal =  set([data[bestPartition] for data in dataset])
78      del(labels[bestPartition])
79      for value in uniqueVal:
80          subLabels = labels[:]
81          branchNode = createTree(splitDataset(dataset, bestPartition, value), subLabels)
82          branchNode.assignDecision(value)
83          newTree.addBranch(branchNode)
84      return newTree
```

# 5  Naive Bayes Algorithm

## 5.1  Basic Method

Given that the input vector $X \subseteq \mathbb{R}^n$, where $X = \begin{bmatrix} X^{(1)} & X^{(2)} & \cdots & X^{(n)} \end{bmatrix}^T$, and the relevant output class label set $Y \in \{c_1, c_2, \ldots, c_K\}$. The joint distribution $P(X, Y)$ generates the data outcome

$$T = \{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\} \tag{9}$$

Applying Bayes rule we construct the algorithm for determining the label for an arbitrary new input.

Assume that the input vector is $x = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(n)} \end{bmatrix}$, the probable label set is $\boldsymbol{c} = \{c_1, c_2, \ldots, c_K\}$, then for $c_i \in \boldsymbol{c}$

$$P(X = x | Y = c_i) = P\left(X^{(1)} = x^{(1)}, X^{(2)} = x^{(2)}, \ldots, X^{(n)} = x^{(n)} | Y = c_i\right) \tag{10}$$

Assume that the variables in $X$ are mutually independant, then the posteriori distribution is

$$P(Y = c_i | X = x) = \frac{P(X = x | Y = c_i)P(Y = c_i)}{\sum_i^k P(X = x | Y = c_i)P(Y = c_i)} = \frac{\prod_j P(X^{(j)} = x^{(j)} | Y = c_i)P(Y = c_i)}{\sum_i \prod_j P(X^{(j)} = x^{(j)} | Y = c_i)P(Y = c_i)}$$

The optimal choice for $c_i$ is

$$
y = f(x) = \arg\max_{c_i} P(Y = c_i | X = x) = \arg\max_{c_i} \frac{\prod_j P(X^{(j)} = x^{(j)} | Y = c_i) P(Y = c_i)}{\sum_i \prod_j P(X^{(j)} = x^{(j)} | Y = c_i) P(Y = c_i)}
$$
$$
= \arg\max_{c_i} \prod_j P(X^{(j)} = x^{(j)} | Y = c_i) P(Y = c_i) \tag{11}
$$

Assume that $x^{(j)} \in \boldsymbol{a_j} = \{a_{j,1}, a_{j,2}, \ldots, a_{j,S_j}\}$, where $\boldsymbol{a_j}$ is the probable value set, then the apriori probability is

$$
P(Y = c_i) = \frac{\sum_{k=1}^{N} I(y_k = c_i)}{N}
$$
$$
P\left(X^{(j)} = a_{j,l} | Y = c_i\right) = \frac{\sum_{k=1}^{N} I(x_k^{(j)} = a_{j,l}, y_k = c_i)}{\sum_{k=1}^{N} I(y_k = c_i)} \quad 1 \le l \le S_j,\ 1 \le j \le n,\ 1 \le k \le K
$$