# Practical 4

**Google colab link:** https://colab.research.google.com/drive/1kx2fQ-NvedFtb05ITteRVziiPhTm8OGM#scrollTo=Xbb2SZUcC3vH&line=3&uniqifier=1

**Aim: Decision Tree - ID3 Algorithm, Naive Bayes Theorem**

**Theory:**

Decision trees are assigned to the information-based learning algorithms which use different measures of information gain for learning it is the most powerful and popular tool for classification and prediction. We can use decision trees for issues where we have continuous but also categorical input and target features. Decision tree is a flowchart like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class labe

Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset. Naive bayes does quite well when the training data doesn't contain all possibilities so it can be very good with low amounts of data. Decision trees work better with lots of data compared to Naive Bayes. Naive Bayes is used a lot in robotics and computer vision, and does quite well with those task

**Code:**

```
[1]  from google.colab import files
     uploaded = files.upload()
```

Choose Files   3-dataset.csv
 • **3-dataset.csv**(n/a) - 414 bytes, last modified: 7/27/2021 - 100% done
Saving 3-dataset.csv to 3-dataset.csv

```
[2]  import pandas as pd
     import math
     import numpy as np

     data = pd.read_csv("3-dataset.csv")
     features = [feat for feat in data]
     features.remove("answer")

     class Node:
         def __init__(self):
             self.children = []
             self.value = ""
             self.isLeaf = False
             self.pred = ""
```

```python
[2]    def entropy(examples):
           pos = 0.0
           neg = 0.0
           for _, row in examples.iterrows():
               if row["answer"] == "yes":
                   pos += 1
               else:
                   neg += 1
           if pos == 0.0 or neg == 0.0:
               return 0.0
           else:
               p = pos / (pos + neg)
               n = neg / (pos + neg)
               return -(p * math.log(p, 2) + n * math.log(n, 2))

       def info_gain(examples, attr):
           uniq = np.unique(examples[attr])
           #print ("\n",uniq)
           gain = entropy(examples)
           #print ("\n",gain)
           for u in uniq:
               subdata = examples[examples[attr] == u]
               #print ("\n",subdata)
               sub_e = entropy(subdata)
```

```python
           gain -= (float(len(subdata)) / float(len(examples))) * sub_e
           #print ("\n",gain)
       return gain

   def ID3(examples, attrs):
       root = Node()

       max_gain = 0
       max_feat = ""
       for feature in attrs:
           #print ("\n",examples)
           gain = info_gain(examples, feature)
           if gain > max_gain:
               max_gain = gain
               max_feat = feature
       root.value = max_feat
       #print ("\nMax feature attr",max_feat)
       uniq = np.unique(examples[max_feat])
       #print ("\n",uniq)
       for u in uniq:
           #print ("\n",u)
           subdata = examples[examples[max_feat] == u]
           #print ("\n",subdata)
           if entropy(subdata) == 0.0:
               newNode = Node()
```

```
                newNode.isLeaf = True
                newNode.value = u
                newNode.pred = np.unique(subdata["answer"])
                root.children.append(newNode)
            else:
                dummyNode = Node()
                dummyNode.value = u
                new_attrs = attrs.copy()
                new_attrs.remove(max_feat)
                child = ID3(subdata, new_attrs)
                dummyNode.children.append(child)
                root.children.append(dummyNode)
    return root

def printTree(root: Node, depth=0):
    for i in range(depth):
        print("\t", end="")
    print(root.value, end="")
    if root.isLeaf:
        print(" -> ", root.pred)
    print()
    for child in root.children:
        printTree(child, depth + 1)

root = ID3(data, features)

printTree(root)
```

**Output:**

```
outlook
        overcast ->  ['yes']

        rain
                wind
                        strong ->  ['no']

                        weak ->  ['yes']

        sunny
                humidity
                        high ->  ['no']

                        normal ->  ['yes']
```