

## Decision Tree Classifier using Information Gain and Entropy Calculations

```
In [904. From ucimlrepo import fetch_ucirepo
import pandas
import math
import random
import numpy

import copy

random.seed(13) # reproducibility

# fetch dataset
mushroom = fetch_ucirepo(id=73)

# data (as pandas dataframes)
X = mushroom.data.features
y = mushroom.data.targets

def subset(arr, feature, value):
    index = list(column_names).index(feature)
    return [item for item in arr if item[index] == value]

def split(arr, base = findroot(), unused = list(column_names), n =0, max_depth=4):
    if n == max_depth:
        return
    feature = select_infogain(arr, unused, base)[0]
    node = [feature, []]
    tree.append(node)
    unique = set(X[feature].values)
    i = unused.index(feature)
    new = unused[:i]+unused[i+1:]

    for value in unique:
        if not pandas.isnull(value):
            sub = subset(arr, feature, value)
            next_f = select_infogain(sub, new, base)[0]
            if new_ent[sub][0] == 0:
                if new_ent[sub][1] == 0:
                    node[i].append([value, 'p'])
                else:
                    node[i].append([value, 'e'])
            else:
                node[i].append([value, next_f, new_ent[sub][1], new_ent[sub][2]])
            split(sub, new_ent[sub][0], unused[:i], n + 1, max_depth)

def new_ent(arr):
    p_count = 0
    e_count = 0
    for item in arr:
        if item[-1] == 'p':
            p_count+=1
        elif (item[-1])=='e':
            e_count+=1
    total = e_count + p_count
    if total == 0 or e_count==0 or p_count ==0:
        base_entropy = 0
    else:
        base_entropy = -1*(((e_count/total)*math.log2(e_count/total)) + ((p_count/total)*math.log2(p_count/total)) )
    return [base_entropy, e_count, p_count]

def findroot(): #finds the root of the tree we will use
    values = y['poisonous'].values
    pcount = 0
    ecount = 0
    for v in values:
        if (v == 'p'):
            pcount+=1
        elif(v == 'e'):
            ecount+=1
    total = ecount + pcount
    base_entropy = -1*((ecount/total)*math.log2(ecount/total)) + ((pcount/total)*math.log2(pcount/total)) )
    return base_entropy
# need to create seth which checks each feature then checks for all unique values
#pass array down with the base entropy, and then check each columns unique values and then

def select_infogain(arr, unused_features, base):
    infogains = []
    for f in unused_features:
        infogain = base
        unique = set(list(X[f].values)))
        for char in unique:
            nprop = n.proportion(arr, f, char)
            char_entropy = entropy(arr, f, char)
            infogain -= (nprop * char_entropy)
        infogains.append(infogain)
    i = infogains.index(max(infogains))
    return [unused_features[i], max(infogains)]

# returns proportion of how many chars eg. a's in the category (column) eg. capszie array, a returns number of a's/total
def n_proportion(arr, column, value):
    index = (list(column_names).index(column))
    count = 0
    for item in arr:
        if item[index] == value:
            count+=1
    if len(arr) ==0:
        return 0
    return count/len(arr)

def entropy(arr, column, value):
    poison = poison_prop(arr, column, value)
    not_poison = 1 - poison
    if poison == 0 or poison == 1:
        return 0
    if not_poison == 0 or not_poison == 1:
        return 0

    entropy = -1*((poison*math.log2(poison))+(not_poison*math.log2(not_poison)))

    return entropy

def poison_prop(arr, column, value): #returns proportion of poisonous when u pass a category and a value, eg. "capsize array", "a"
    index = (list(column_names).index(column))
    value_count = 0
    poison_count = 0

    for item in arr:
        if item[index] == value:
            value_count+=1
            if item[-1]=='p':
                poison_count+=1
    if value_count ==0:
        return 0
    return poison_count/value_count

def print_tree(t, depth=1):
    d = depth
    if isinstance(t, list)==True:
        for i in range(len(t)):
            item = t[i]
            if len(item)>2:
                item = item[0:2]
                if isinstance(item[1],list)==True:
                    print(" " + str(item[0]))
                    print_tree(item[1], d+1)
                else:
                    if item[1] == 'p':
                        print("[" + "d + item[0]" + " -> " + "POISONOUS")
                    elif item[1] == 'e':
                        print("[" + "d + item[0]" + " -> " + "EDIBLE")
                    else:
                        print("[" + "d + item[0]" + " -> " + item[1])

column_names = X.columns.to_numpy()
classified_data = []
index = (list(column_names).index(column))
for i in range(len(X[feature_columns])):
    data = []
    for j in range (len(column_names)):
        data.append(X[column_names[j]].values[i])
    found = False
    for d in data:
        if pandas.isnull(d):
            found = True
        if found == False:
            data.append(y[poisonous].values[i])
            classified_data.append(data)

# this puts the data in lists inside a list, classified data, with each mushroom's features in order.
```

```
In [905. #Output for treedepth 2
tree = []
split(classified_data, findroot(), list(column_names), 0, 2)
print_tree(tree)
```

```
odor
| | y --> POISONOUS
| | l --> EDIBLE
| | a --> EDIBLE
| | c --> POISONOUS
| | n --> spore-print-color
| | p --> POISONOUS
| | s --> POISONOUS
| | m --> POISONOUS
| | f --> POISONOUS
| | k --> EDIBLE
| | b --> POISONOUS
| | y --> POISONOUS
| | w --> cap-color
| | n --> EDIBLE
| | o --> POISONOUS
| | b --> POISONOUS
| | r --> POISONOUS
| | u --> POISONOUS
```

```
In [906. #Output for treedepth 3
tree = []
split(classified_data, findroot(), list(column_names), 0, 3)
print_tree(tree)
```

```
odor
| | y --> POISONOUS
| | l --> EDIBLE
| | a --> EDIBLE
| | c --> POISONOUS
| | n --> spore-print-color
| | p --> POISONOUS
| | s --> POISONOUS
| | m --> POISONOUS
| | f --> POISONOUS
| | k --> EDIBLE
| | b --> POISONOUS
| | y --> POISONOUS
| | w --> cap-color
| | n --> EDIBLE
| | o --> POISONOUS
| | b --> POISONOUS
| | r --> POISONOUS
| | u --> POISONOUS
```

```
In [907. #Output for treedepth 4
tree = []
split(classified_data, findroot(), list(column_names), 0, 4)
print_tree(tree)
```

```
odor
| | y --> POISONOUS
| | l --> EDIBLE
| | a --> EDIBLE
| | c --> POISONOUS
| | n --> spore-print-color
| | p --> POISONOUS
| | s --> POISONOUS
| | m --> POISONOUS
| | f --> POISONOUS
| | k --> EDIBLE
| | b --> POISONOUS
| | y --> POISONOUS
| | w --> cap-color
| | n --> EDIBLE
| | o --> POISONOUS
| | b --> POISONOUS
| | r --> POISONOUS
| | u --> POISONOUS
```

Code for Question 1 (a) and (b) where decision tree is made using the dataset provided and tested for treedepths 2,3,4 For the code I printed the tree as a list of lists where the first element in the list defines the feature at that node and the following elements in the list show the outcome of the different treedepths.

Implemented decision tree with infogain splitting criteria, removed nan values and used random.seed(10) for questions (c) and (d) implementation. Simple tree print without indents. First split: Odor Second split: Spore-print-color Third split: Cap-color

The tree representation is not indented as the tree diagram was not required either way. However to read the tree, first look at odor and only if odor = 'n' then look at spore-print-color and so on. Full accuracy when tree depth >=3

```
In [909. def newdata(): # makes new data for one mushroom in form same as classified data, features are randomised, can be anything.
    new_mushroom = []

    for feature in column_names:
        u = list(set(list(X[feature].values)))
        unique = []
        for value in u:
            if not pandas.isnull(value):
                unique.append(value)
        i = random.randrange(0,len(unique))
        new_mushroom.append(unique[i])

    return new_mushroom
```

```
def testing(mushroom, full_tree,treedepth):
    count = 1
    for item in full_tree:
        index = (list(column_names).index(item[0]))
        value = mushroom[index]
        items = item[1]
        for prediction in items:
            if value == prediction[0]:
                if len(prediction) ==2:
                    if prediction[1] == 'p':
                        return "POISONOUS"
                    else
                        return "EDIBLE"
            elif treedepth == count:
                if prediction[2] > prediction[3]:
                    return "EDIBLE"
                else:
                    return "POISONOUS"
            count+=1
```

```
In [910. #testing
```

```
"""
TESTING FOR NEW RANDOM DATA use this code instead of the values I provided below in lists (i,l,m,n)
i = newdata()
l = newdata()
a = newdata()
t = newdata()
"""

"""
run for visuals
print(column_names)
tree=[]
split(classified_data, findroot(), list(column_names), 0, 4)
print_tree(tree)
print(i)
print(l)
print(m)
print(t)
"""

#test for specific data so that the results are reproducible.
i = ['e', 'y', 'u', 'f', 'p', 'a', 'c', 'n', 'b', 'e', 'e', 'f', 'g', 'f', 'g', 'p', 'p', 'w', 'n', 'e', 'u', 'y', 'l']
l = ['e', 'y', 'p', 't', 'n', 'f', 'w', 'n', 'e', 'e', 'b', 'k', 'f', 'o', 'b', 'p', 'o', 'o', 'l', 'n', 'y', 'p']
a = ['f', 'f', 'f', 'f', 'w', 'n', 'b', 'e', 'b', 'f', 'k', 'p', 'w', 'p', 'w', 't', 'n', 'b', 'e', 'g']
t = ['k', 'y', 'u', 'f', 'n', 'f', 'w', 'n', 'o', 't', 'b', 'y', 'k', 'b', 'p', 'p', 'o', 't', 'e', 'o', 'w', 'l']

tree = []
split(classified_data, findroot(), list(column_names), 0, 4)
print(testing(i,tree,4))

tree = []
split(classified_data, findroot(), list(column_names), 0, 2)
print(testing(m,tree,2))

tree = []
split(classified_data, findroot(), list(column_names), 0, 3)
print(testing(l,tree,3))

tree = []
split(classified_data, findroot(), list(column_names), 0, 1)
print(testing(t,tree,1))

POISONOUS
POISONOUS
EDIBLE
EDIBLE
```

Question 1(c): Testing

I tested my algorithm of the decision tree by firstly making a method which constructs a new mushroom with data in the same formatting as my classified\_data variable ( ['a','k','...'] etc. where each index corresponds to a column eg. 'a' corresponds to cap-size).

I randomised this using random.randrange and random.seed(10) to ensure reproducibility. With this new randomised mushroom I ran the method testing(mushroom, treedepth) which uses the tree made by the full dataset to make a prediction on a new mushroom.

This is done by traversing through the tree until either reaching a leaf node or making a prediction based on the proportion of poisonous mushrooms and edible mushrooms in that level of the tree as to not exceed treedepth. When treedepth >= 4 there is a 100% accuracy as level 4 contains solely of leaf nodes.

Random.seed(13) for reproducibility. However, code can be tested further by making more data. Random seed is set at the beginning of the main code (q1 a&b where libraries are imported and seed is initialised.

Make sure to set tree = [] if testing to get correct results.

for the sake of reproducibility i have used specific data, but using my newdata() function.

```
In [912. def evaluation(full_set, full_tree, treedepth):
    correct = 0
    incorrect = 0
    total = 0
```

```
    for data in full_set:
        status = testing(data, full_tree, treedepth)
        if (data[-1] == 'p') and (status == "POISONOUS"):
            correct+=1
        elif (data[-1] == 'e') and (status == "EDIBLE"):
            correct+=1
        total+=1
    return correct/total
```

```
In [913. full = []
```

```
    for k in classified_data:
        full.append(k)

    testing_data = []

    for i in range((3000)):
        random.i = random.randrange(0, len(full))
        testing_data.append(full[random.i])
        full.pop(random.i)

    for i in range(1,5):
        tree = []
        split(full, findroot(), list(column_names), 0, i)
        accuracy_i = evaluation(testing_data, tree, i)
        print("The testing accuracy for treedepth " + str(i) + " is " + str(accuracy_i))

The testing accuracy for treedepth 1 is 0.9833333333333333
The testing accuracy for treedepth 2 is 0.9966666666666667
The testing accuracy for treedepth 3 is 1.0
The testing accuracy for treedepth 4 is 1.0
```

Full accuracy when tree\_depth >=3 because I removed the missing values from the dataset. Had I not done this the accuracy would only be 100% when tree\_depth==4.

I implemented a function which evaluates the whole decision tree procedure. First it forms two sets based on the full data set randomly one being the testing\_data set and one being the training data set (full). It uses the training data set to form a decision tree that is split based on information gain and entropy.

Once the tree is formed we can test the accuracy of the data by using the test function I made in Q1(C) on each list of data. This will give us a prediction on whether the mushroom is poisonous or not and we can compare this to the actual data of the mushroom to see if the prediction was correct. Doing this for many mushrooms/data gives us an accuracy score. (number of prediction correct/total predictions).

For treedepth 1 the accuracy is 0.98425 For treedepth 2 the accuracy is 0.99725 For treedepth 4 it is 1.0 and if treedepth >=3 then the accuracy will be 1.0 (100%).

There are several other ways the decision tree can be implemented using different splitting criteria. A common splitting criteria is the gini impurity splitting criteria. Altering the splitting criteria will change how the decision tree determines which feature to split into at what level of the tree, thus resulting in different partitions and a completely different tree and thus with a different accuracy.

Gini impurity is measured as a number from 0 to 1 where 0 indicates all elements are part of the same class and 1 indicating all elements are randomly amongst a variety of classes. 0.5 shows a normal distribution of the elements amongst the classes.

The gini impurity can be calculated by calculating gini coefficients for the subnodes of a node and then calculating the impurity of each node by calculating the weighted average of the distribution amongst the subnodes in the non-binary tree. The idea here is to minimise the number of incorrect nodes when classifying the mushrooms.

The gini impurity splitting criteria is based around the probability of data being misclassified. The entropy, information gain splitting criteria both focus on the impurity in a dataset. Gini entropy has a more favourable computing power meaning it is preferable when datasets are larger as the entropy splitting criteria will be more inefficient. However, using entropy typically takes longer but is also more accurate as it seeks the maximum information gain and the least entropy and minimises inaccuracy. For larger datasets, gini impurity is favourable due to the computing power. Additionally gini impurity is favourable for datasets where the classes are relatively balanced, but using the informative gain is favourable when classes are more randomly distributed!

Using gini impurity and entropy will usually give two different decision trees, which is important to note. This is because the splitting criteria are different which means the nodes selected are likely different. This means the decision tree will usually have a different accuracy depending on whether you use entropy/information gain or gini impurity to implement the algorithm.

```
In [899. training_data = []
for k in classified_data:
    training_data.append(k)

for i in range(1,5):
    tree = []
    split(training_data, findroot(), list(column_names), 0, i)
    accuracy_i = evaluation(training_data, tree, i)
    print("The training accuracy for treedepth " + str(i) + " is " + str(accuracy_i))
```

The training accuracy for treedepth 1 is 0.984408221197732
The training accuracy for treedepth 2 is 0.99716513112686
The training accuracy for treedepth 3 is 1.0
The training accuracy for treedepth 4 is 1.0

My evaluation can explain whether my tree is overfitting by now testing the training accuracy.

The training accuracy for: treedepth 1 is 0.984408221197732 The training accuracy for treedepth 2 is 0.99716513112686 The training accuracy for treedepth 3 is 1.0 The training accuracy for treedepth 4 is 1.0

And we know from before: The testing accuracy for treedepth 1 is 0.9833333333333333 The testing accuracy for treedepth 2 is 0.9966666666666667 The testing accuracy for treedepth 3 is 1.0 The testing accuracy for treedepth 4 is 1.0

Underfitting data is typically when training accuracy is low and it is usually overfitting when testing accuracy is low. In this case both the training accuracy and testing accuracy are high which means it is neither underfitting nor overfitting!