

DATA ANALYTICS ASSIGNMENT

NIVEDITHA C U

RACHANA H S

PES1201701640

PES1201701726

SECTION 'C'

SECTION 'E'

PROBLEM STATEMENT:

BUILDING DECISION TREE FOR A DATASET

SCREENSHOTS

```
In [1]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

import random
from pprint import pprint
```

```
In [2]: df = pd.read_csv("iris_data.csv")
df = df.drop("Id", axis=1)
df = df.rename(columns={"Species": "label"})
```

```
In [3]: df.head()
```

```
Out[3]:
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	label
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

```
In [4]: def train_test_split(df, test_size):

    if isinstance(test_size, float):
        test_size = round(test_size * len(df))

    indices = df.index.tolist()
    test_indices = random.sample(population=indices, k=test_size)

    test_df = df.loc[test_indices]
    train_df = df.drop(test_indices)

    return train_df, test_df
```

```
In [5]: random.seed(0)
train_df, test_df = train_test_split(df, test_size=20)
```

```
In [6]: data = train_df.values
data[:5]
```

```
Out[6]: array([[5.1, 3.5, 1.4, 0.2, 'Iris-setosa'],
 [4.9, 3.0, 1.4, 0.2, 'Iris-setosa'],
 [4.7, 3.2, 1.3, 0.2, 'Iris-setosa'],
 [4.6, 3.1, 1.5, 0.2, 'Iris-setosa'],
 [5.0, 3.6, 1.4, 0.2, 'Iris-setosa']], dtype=object)
```

```
In [7]: def check_purity(data):

    label_column = data[:, -1]
    unique_classes = np.unique(label_column)

    if len(unique_classes) == 1:
        return True
    else:
        return False
```

```
In [8]: def classify_data(data):

    label_column = data[:, -1]
    unique_classes, counts_unique_classes = np.unique(label_column, return_counts=True)

    index = counts_unique_classes.argmax()
    classification = unique_classes[index]

    return classification
```

```
In [9]: def get_potential_splits(data):

    potential_splits = {}
    _, n_columns = data.shape
    for column_index in range(n_columns - 1):      # excluding the last column which is the label
        potential_splits[column_index] = []
        values = data[:, column_index]
        unique_values = np.unique(values)

        for index in range(len(unique_values)):
            if index != 0:
                current_value = unique_values[index]
                previous_value = unique_values[index - 1]
                potential_split = (current_value + previous_value) / 2

            potential_splits[column_index].append(potential_split)

    return potential_splits
```

```
In [10]: def split_data(data, split_column, split_value):

    split_column_values = data[:, split_column]

    data_below = data[split_column_values <= split_value]
    data_above = data[split_column_values > split_value]

    return data_below, data_above
```

```
In [11]: def calculate_entropy(data):

    label_column = data[:, -1]
    _, counts = np.unique(label_column, return_counts=True)

    probabilities = counts / counts.sum()
    entropy = sum(probabilities * -np.log2(probabilities))

    return entropy
```

In [12]: `def calculate_overall_entropy(data_below, data_above):`

```
    n = len(data_below) + len(data_above)
    p_data_below = len(data_below) / n
    p_data_above = len(data_above) / n

    overall_entropy = (p_data_below * calculate_entropy(data_below)
                       + p_data_above * calculate_entropy(data_above))

    return overall_entropy
```

In [13]: `def determine_best_split(data, potential_splits):`

```
    overall_entropy = 9999
    for column_index in potential_splits:
        for value in potential_splits[column_index]:
            data_below, data_above = split_data(data, split_column=column_index, split_value=value)
            current_overall_entropy = calculate_overall_entropy(data_below, data_above)

            if current_overall_entropy <= overall_entropy:
                overall_entropy = current_overall_entropy
                best_split_column = column_index
                best_split_value = value

    return best_split_column, best_split_value
```

In [14]: `example_tree = {"petal_width <= 0.8": ["Iris-setosa", { "petal_width <= 1.65": [{"petal_length <= 4.9": ["Iris-versicolor", "Iris-virginica"], "petal_length > 4.9": ["Iris-setosa"]}]}]`

```
In [16]: tree = decision_tree_algorithm(train_df, max_depth=3)
pprint(tree)

{'PetalWidthCm <= 0.8': ['Iris-setosa',
                        {'PetalWidthCm <= 1.65': [{'PetalLengthCm <= 4.95': ['Iris-versicolor',
                                                                              'Iris-virginica']],
                        'Iris-virginica']}]}
```

```
In [17]: #classification

example = test_df.iloc[0]
example
```

```
Out[17]: SepalLengthCm      5.1
SepalWidthCm      2.5
PetalLengthCm      3
PetalWidthCm      1.1
label      Iris-versicolor
Name: 98, dtype: object
```

```
In [18]: def classify_example(example, tree):
question = list(tree.keys())[0]
feature_name, comparison_operator, value = question.split(" ")

# ask question
if example[feature_name] <= float(value):
    answer = tree[question][0]
else:
    answer = tree[question][1]

# base case
if not isinstance(answer, dict):
    return answer

# recursive part
else:
    residual_tree = answer
    return classify_example(example, residual_tree)
```

```
In [19]: classify_example(example, tree)
```

```
Out[19]: 'Iris-versicolor'
```

```

In [15]: def decision_tree_algorithm(df, counter=0, min_samples=2, max_depth=5):

    # data preparations
    if counter == 0:
        global COLUMN_HEADERS
        COLUMN_HEADERS = df.columns
        data = df.values
    else:
        data = df

    # base cases
    if (check_purity(data)) or (len(data) < min_samples) or (counter == max_depth):
        classification = classify_data(data)

        return classification

    # recursive part
    else:
        counter += 1

        # helper functions
        potential_splits = get_potential_splits(data)
        split_column, split_value = determine_best_split(data, potential_splits)
        data_below, data_above = split_data(data, split_column, split_value)

        # instantiate sub-tree
        feature_name = COLUMN_HEADERS[split_column]
        question = "{} <= {}".format(feature_name, split_value)
        sub_tree = {question: []}

        # find answers (recursion)
        yes_answer = decision_tree_algorithm(data_below, counter, min_samples, max_depth)
        no_answer = decision_tree_algorithm(data_above, counter, min_samples, max_depth)

        # If the answers are the same, then there is no point in asking the question.
        # This could happen when the data is classified even though it is not pure
        # yet (min_samples or max_depth base case).
        if yes_answer == no_answer:
            sub_tree[question] = yes_answer
        else:
            sub_tree[question].append(yes_answer)
            sub_tree[question].append(no_answer)

    return sub_tree

```

```

In [20]: #accuracy

def calculate_accuracy(df, tree):

    df["classification"] = df.apply(classify_example, axis=1, args=(tree,))
    df["classification_correct"] = df["classification"] == df["label"]

    accuracy = df["classification_correct"].mean()

    return accuracy

```

```

In [21]: accuracy = calculate_accuracy(test_df, tree)
accuracy

```

```

Out[21]: 0.95

```