**Machine Learning & Data-Mining : Research Project**

**Implementation and experiments**

**Predictions of diseases using decision trees**

**References for this project :**
- Chronic kidney disease diagnosis usin decision trees algorithms (2021)
(https://bmcnephrol.biomedcentral.com/articles/10.1186/s12882-021-02474-z)
- Simple Prediction of Type 2 Diabetes Mellitus via Decision Tree
Modeling (2017) (https://brieflands.com/articles/ircrj-10657.pdf)
- Decision tree model in the diagnosis of breast cancer (2017)
(https://ieeexplore.ieee.org/document/8789297)
- Early Prediction of Heart Disease Using
Decision Tree Algorithm (2017) (https://www.researchgate.net/profile/Safish-Mary/publication/315023624_Early_Prediction_of_Heart_Disease_Using_Decision_Tree_Algorithm/links/58c84b57aca2723ab16eba60/Early-Prediction-of-Heart-Disease-Using-Decision-Tree-Algorithm.pdf)

I will get help with the code of the first assignment on decision trees for my implementation.

I am re-implementing what I coded with another dataset. I used the dataset « Breast cancer Wisconsin » : https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data

Let's highlight the key components of the code and their functionalities:

- Libraries :

```
6    import pandas as pd
7    from typing import Union
8    import matplotlib.pyplot as plt
9    import numpy as np
10   from sklearn.tree import DecisionTreeClassifier, plot_tree
11   from sklearn.metrics import accuracy_score
```

- Data Loading and Preprocessing:

The load_data function loads the Breast Cancer dataset from a CSV file. This function drops unnecessary columns from the dataset and features and targets are extracted.

```
201   def load_data():
202       BreastCancerDataset = pd.read_csv('BreastCancerDataset.csv').drop(['Unnamed: 32'], axis=1)
203       BreastCancerFeatures = BreastCancerDataset.drop(['id', 'diagnosis'], axis=1).values
204       BreastCancerTargets = BreastCancerDataset['diagnosis'].values
205       BreastCancerClasses = BreastCancerDataset['diagnosis'].unique()
206       return BreastCancerFeatures, BreastCancerTargets, BreastCancerClasses
```

- Prior probability function

The prior function calculates the prior probability of each class type in a classification problem. It represents the distribution of classes in the dataset before considering any features.

```
35    def prior(targets: np.ndarray, classes: list) -> np.ndarray:
36        '''
37        Calculate the prior probability of each class type
38        given a list of all targets and all class types
39        '''
40        class_counts = np.zeros(len(classes))
41
42        for target in targets:
43            class_counts[target] += 1
44
45        total_samples = len(targets)
46        class_probabilities = class_counts / total_samples
47
48        return class_probabilities
```

- Splitting Data:

The split_train_test function shuffles and splits the dataset into training and testing sets based on a specified ratio. I am using the same split_train_test function from the tools of the assignment.

```
13    def split_train_test(features: np.ndarray, targets: np.ndarray,
14        train_ratio:float=0.8) -> Union[tuple, tuple]:
15        '''
16        Shuffle the features and targets in unison and return
17        two tuples of datasets, first being the training set,
18        where the number of items in the training set is according
19        to the given train_ratio
20        '''
21        np.random.seed(888)
22        p = np.random.permutation(features.shape[0])
23        features = features[p]
24        targets = targets[p]
25
26        split_index = int(features.shape[0] * train_ratio)
27
28        train_features, train_targets = features[0:split_index, :],\
29        targets[0:split_index]
30        test_features, test_targets = features[split_index:-1, :],\
31            targets[split_index: -1]
32
33        return (train_features, train_targets), (test_features, test_targets)
34
```

- Decision Tree Functions:
  ➢ split_data Function:

This function splits the dataset and targets into two separate datasets based on a given feature and threshold.

```
51    def split_data(
52        features: np.ndarray,
53        targets: np.ndarray,
54        split_feature_index: int,
55        theta: float
56    ) -> Union[tuple, tuple]:
57        '''
58        Split a dataset and targets into two seperate datasets
59        where data with split_feature < theta goes to 1 otherwise 2
60        '''
61        split1 = features[:, split_feature_index] < theta
62        split2 = features[:, split_feature_index] >= theta
63
64        features_1 = features[split1]
65        targets_1 = targets[split1]
66
67        features_2 = features[split2]
68        targets_2 = targets[split2]
69
70        return (features_1, targets_1), (features_2, targets_2)
```

➢ gini_impurity Function:

This function calculates the Gini impurity for a set of targets.

```python
73    def gini_impurity(targets: np.ndarray, classes: list) -> float:
74        '''
75        Calculate:
76            i(S_k) = 1/2 * (1 - sum_i P{C_i}**2)
77        '''
78
79        square = prior(targets, classes)
80        for n in range(len(square)): square[n] = square[n]**2
81
82        return 1/2 * (1 - sum(square))
```

➢ weighted_impurity Function:

This function calculates the weighted sum of Gini impurities for two branches.

```python
85    def weighted_impurity(
86        t1: np.ndarray,
87        t2: np.ndarray,
88        classes: list
89    ) -> float:
90        '''
91        Given targets of two branches, return the weighted
92        sum of gini branch impurities
93        '''
94        g1 = gini_impurity(t1,classes)
95        g2 = gini_impurity(t2,classes)
96        n = t1.shape[0] + t2.shape[0]
97
98        i = (t1.shape[0] * g1 + t2.shape[0] * g2) / n
99        return i
```

➢ total_gini_impurity Function:

This function calculates the total Gini impurity which returns the weighted impurity given the dataset and threshold to split on.

```python
102    def total_gini_impurity(
103        features: np.ndarray,
104        targets: np.ndarray,
105        classes: list,
106        split_feature_index: int,
107        theta: float
108    ) -> float:
109        '''
110        Calculate the gini impurity for a split on split_feature_index
111        for a given dataset of features and targets.
112        '''
113
114        (features_1, targets_1),(features_2, targets_2) = split_data(features, targets, split_feature_index, theta)
115        weight = weighted_impurity(targets_1, targets_2, classes)
116
117        return weight
```

➢ brute_best_split Function:

This function finds the best split for the given data by iterating over feature dimensions and thresholds.

```python
120    def brute_best_split(
121        features: np.ndarray,
122        targets: np.ndarray,
123        classes: list,
124        num_tries: int
125    ) -> Union[float, int, float]:
126        '''
127        Find the best split for the given data. Test splitting
128        on each feature dimension num_tries times.
129
130        Return the lowest gini impurity, the feature dimension and
131        the threshold
132        '''
133
134        best_gini, best_dim, best_theta = float("inf"), None, None
135
136        # iterate feature dimensions
137        for i in range(features.shape[1]):
138            features_i = features[:,i]
139
140            # create the thresholds
141            thetas = np.linspace(features_i.min(), features_i.max(), num_tries+2)[1:-1]
142
143            # iterate thresholds
144            for theta in thetas:
145                gini = total_gini_impurity(features, targets, classes, i, theta)
146                if best_gini > gini:
147                    best_gini = gini
148                    best_dim = i
149                    best_theta = theta
150
151        return best_gini, best_dim, best_theta
```

- Decision Tree Trainer Class:

This class handles the training and evaluation of the Decision Tree model of the breast cancer winsconsin dataset. It utilizes the DecisionTreeClassifier from scikit-learn and provides methods for training, accuracy calculation, plotting, making predictions, and generating a confusion matrix.
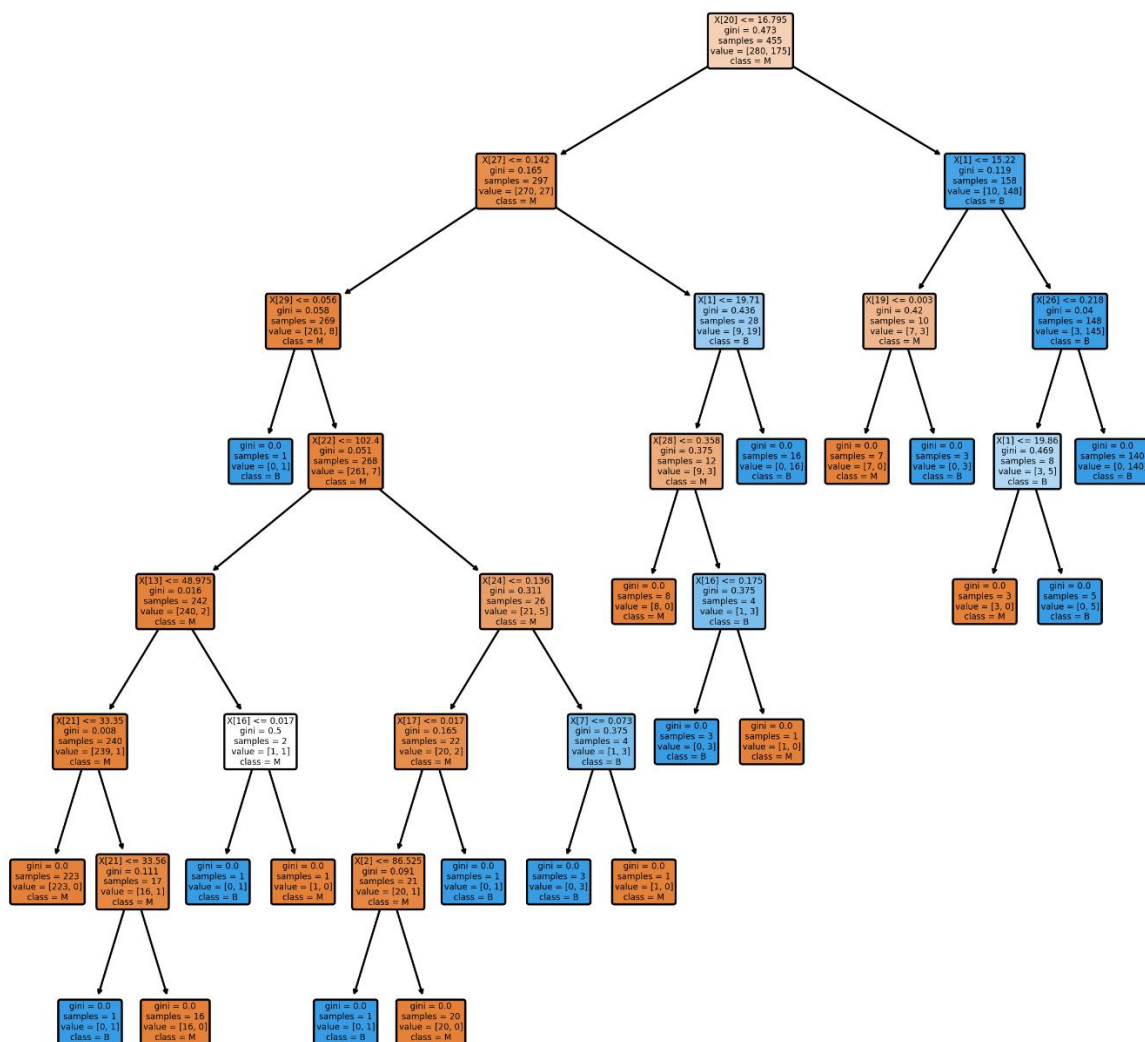
```python
154    class BreastCancerTreeTrainer:
155        def __init__(
156            self,
157            features: np.ndarray,
158            targets: np.ndarray,
159            classes: list = ['M','B'],
160            train_ratio: float = 0.8
161        ):
162            '''
163            train_ratio: The ratio of the Breast Cancer dataset that will
164            be dedicated to training.
165            '''
166            (self.train_features, self.train_targets),\
167                (self.test_features, self.test_targets) =\
168                split_train_test(features, targets, train_ratio)
169
170            self.classes = classes
171            self.tree = DecisionTreeClassifier()
172
173        def train(self):
174            return self.tree.fit(self.train_features, self.train_targets)
175
176        def accuracy(self):
177            return accuracy_score(self.test_targets, self.tree.predict(self.test_features))
178
179        def plot(self):
180            plt.figure(figsize=(10,10), dpi=300)
181            plot_tree(self.tree, filled=True, class_names=self.classes, rounded=True)
182            plt.show()
183
184
185        def guess(self):
186            return self.tree.predict(self.test_features)
187
188        def confusion_matrix(self):
189            predictions = self.guess()
190            num_classes = len(self.classes)
191            cm = np.zeros((num_classes, num_classes), dtype=int)
192
193            for true_label, pred_label in zip(self.test_targets, predictions):
194                cm[true_label][pred_label] += 1
195
196            return cm
```

- Main Script:

In the main, I load data using the load_data function. I instantiate a BreastCancerTreeTrainer object and train a Decision Tree model in order to do prediction and classify if the patient's tumor is benigne or malignant. If it is malignant, it means that the patient has breast cancer.

For the moment, I have this decision tree while taking into consideration all the features of the Breast Cancer Dataset:



The accuracy achieved by my implementation is 95.58%, slightly surpassing the 94.3% reported in the referenced article (https://ieeexplore.ieee.org/document/9442043/). Despite this high accuracy, it's crucial to acknowledge the presence of bugs in my code, potentially leading to inaccurate results.

One identified issue lies in the prior function, which requires modification. The breast cancer dataset's classes, denoted as 'B' for benign and 'M' for malignant, are of string type. This contrasts with the integer type classes in the iris dataset. Consequently, the prior function needs adjustments to handle string-type classes.
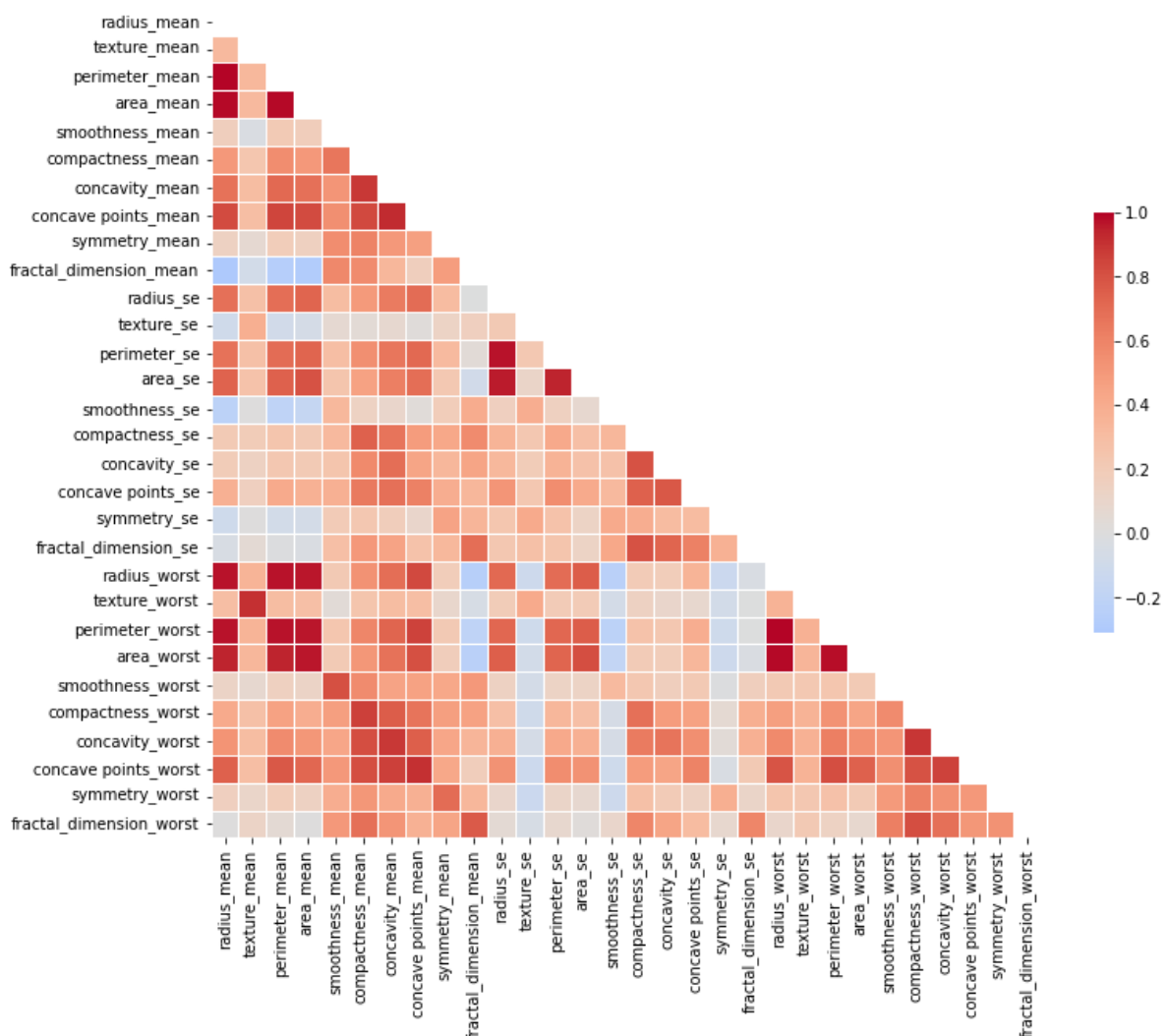
Additionally, the confusion matrix is currently non-functional. Addressing this issue is essential for a comprehensive evaluation of the model's performance.

Aside of this first plot, I tried to see what features of the dataset have the biggest impact on results in order to reduce the number of features (more than 30). By retaining only features that are not correlated, the model can improve in interpretability.
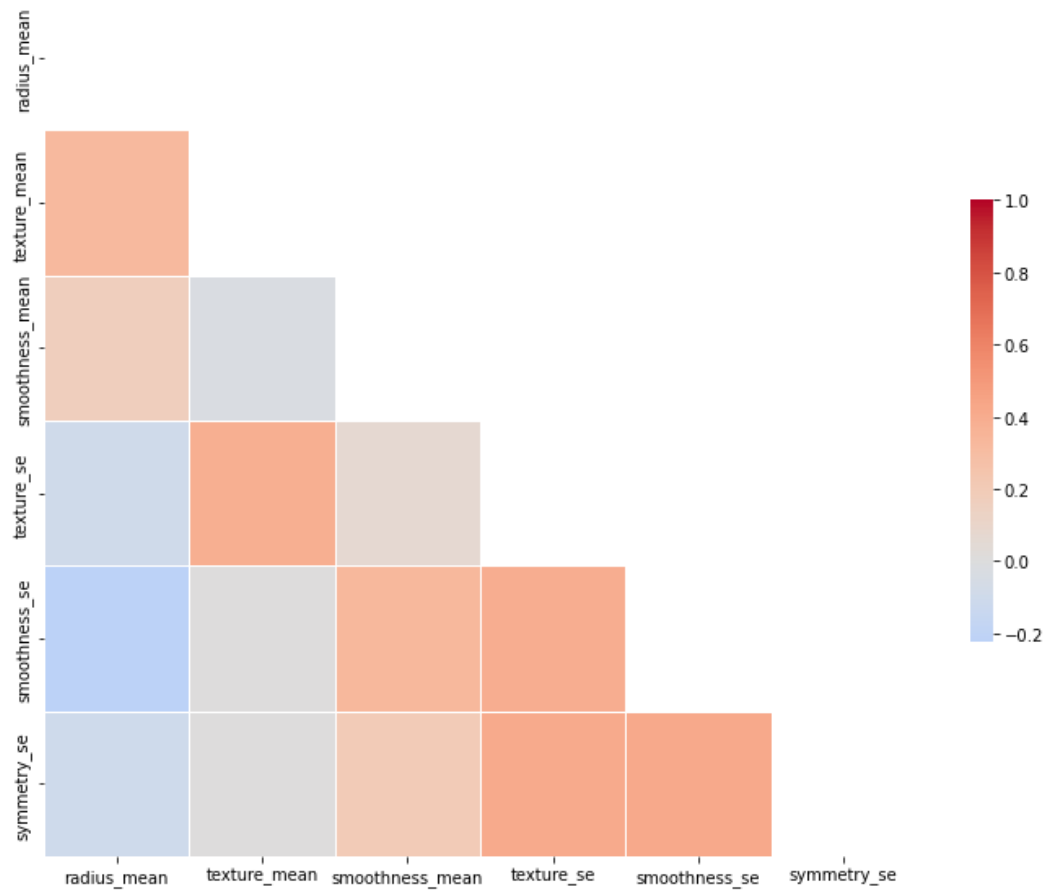
This function helped me in obtaining a correlation map of the Breast Cancer dataset's features.

```python
import seaborn as sns

def feature_correlation_map(dataframe):
    # Calculate the correlation matrix
    corr_matrix = dataframe.corr()

    # Create a mask for the upper triangle
    mask = np.triu(np.ones_like(corr_matrix, dtype=bool))

    # Set up the matplotlib figure
    plt.figure(figsize=(12, 10))

    # Draw the heatmap with the mask
    sns.heatmap(corr_matrix, mask=mask, cmap='coolwarm', vmax=1, center=0,
                square=True, linewidths=.5, cbar_kws={"shrink": 0.5})

    plt.show()
```

I applied the feature_correlation_map to my Breast Cancer Dataset and I obtained this map :



After retaining only the features that has a correlation higher than the threshold of 0.5, we obtain a new correlation map :

**Radius_mean, texture_mean, smoothness_mean, texture_se, smoothness_se and symmetry_se** are the features which are not correlated.

I tried to plot another decision tree but with the features « filtered » :



The accuracy of this decision tree is 88.4%

The next step of my implementation will be to correct bugs that I mention earlier and to understand more deeply how to interpret each decision tree's plot.