

Problem 1

3. It doesn't matter which is used

Problem 2

2. D_2

Problem 3

3. $f(x) = 3$

Problem 4

1. Model 1

M5-L1 Problem 1

In this problem, you will implement a function to calculate gini impurity on an arbitrary input vector.

For reference, the formula for Gini impurity is:

$$\text{Gini}(D) = 1 - \sum_{i=1}^k p_i^2$$

where D is the dataset containing samples from k classes and p_i is the probability of a data point belonging to class i .

Gini Impurity Function

Complete the function `gini(D)` below. It should take as input a 1-D array, where is the number of samples corresponding to each output class.

For example, consider the input array `D = np.array([4, 9, 7, 0, 3])` In this example, there are 5 input classes and 23 total samples. For this input, your function should return 0.707.

Your function should work regardless of the length of the input vector.

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

def gini(D):
    # YOUR CODE GOES HERE
    sum = np.sum(D)
    gini_ = 1 - np.sum(np.square(D / sum))
    return gini_

D = np.array([4, 9, 7, 0, 3])
g = gini(D)
print(f"gini([4,9,7,0,3]) = {g:.3f} (should be about {0.707})")
```

`gini([4,9,7,0,3]) = 0.707 (should be about 0.707)`

More test cases

Compute and print the gini impurity for `D1`, `D2`, `D3`, and `D4`, defined below:

```
In [2]: D1 = np.array([1,0,0])
D2 = np.array([0,0,4])
D3 = np.array([0, 20, 0, 0, 0, 3])
D4 = np.array([6, 6, 6, 6])

for D in [D1, D2, D3, D4]:
    # YOUR CODE GOES HERE
    g = gini(D)
    print(f"gini({D}) = {g:.3f}")
```

```
gini([1 0 0]) = 0.000
gini([0 0 4]) = 0.000
gini([ 0 20 0 0 0 3]) = 0.227
gini([6 6 6 6]) = 0.750
```

Processing math: 100%

M5-L1 Problem 2

Now we will provide a 2D classification dataset and you will learn to use sklearn's decision tree classifier on the data.

First, run the following cell to load the data and import decision tree tools.

- Input: X , size 80×2
- Output: y , size 80

```

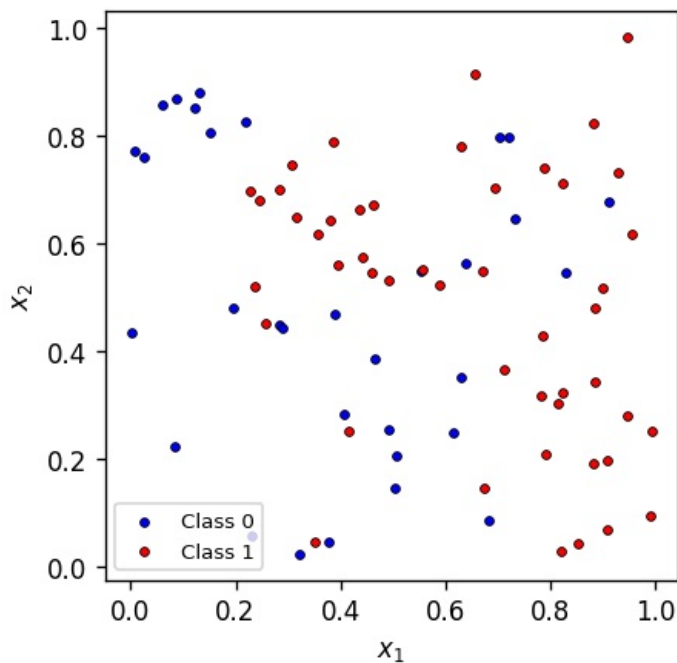
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeClassifier, plot_tree
from matplotlib.colors import ListedColormap

y = np.array([1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0])
x1 = np.array([6.73834679e-01, 3.57095269e-01, 4.42510505e-01, 8.48412660e-02, 2.17890220e-01, 4.60241400e-01,
x2 = np.array([0.14784469, 0.61647661, 0.57595235, 0.2232836, 0.82559199, 0.54569237, 0.73986085, 0.79782627,
X = np.vstack([x1, x2]).T

def plot_data(X,y):
    colors=["blue", "red"]
    for i in range(2):
        plt.scatter(X[y==i,0],X[y==i,1],s=12,c=colors[i],edgecolors="black",linewidths=.5,label=f"Class {i}")
        plt.xlabel("$x_1$")
        plt.ylabel("$x_2$")
        plt.legend(loc="lower left",prop={'size':8})

plt.figure(figsize=(4,4),dpi=120)
plot_data(X,y)
plt.show()

```



Create and fit a decision tree classifier

Create an instance of a `DecisionTreeClassifier()` with `max_depth` of 5. Fit this to the data `X`, `y`.

For more details, consult: <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

```
In [2]: # YOUR CODE GOES HERE
model = DecisionTreeClassifier(max_depth = 5)
model.fit(X, y)
```

```
Out[2]: DecisionTreeClassifier
DecisionTreeClassifier(max_depth=5)
```

Making new predictions using your model

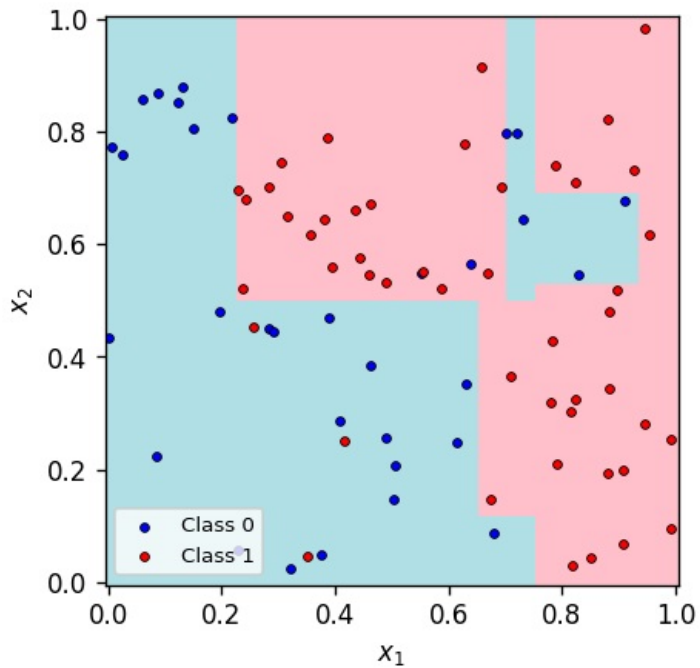
Now use the decision tree you trained to evaluate on the meshgrid of points `X_test` as indicated below. The code here will generate a plot showing the decision boundaries created by the model.

```
In [3]: vals = np.linspace(0,1,100)
x1grid, x2grid = np.meshgrid(vals, vals)

X_test = np.vstack([x1grid.flatten(), x2grid.flatten()]).T

# YOUR CODE GOES HERE
# compute a prediction, `pred` for the input `X_test`
pred = model.predict(X_test)

plt.figure(figsize=(4,4),dpi=120)
bgcolors = ListedColormap(["powderblue","pink"])
plt.pcolormesh(x1grid, x2grid, pred.reshape(x1grid.shape), shading="nearest", cmap=bgcolors)
plot_data(X,y)
plt.show()
```



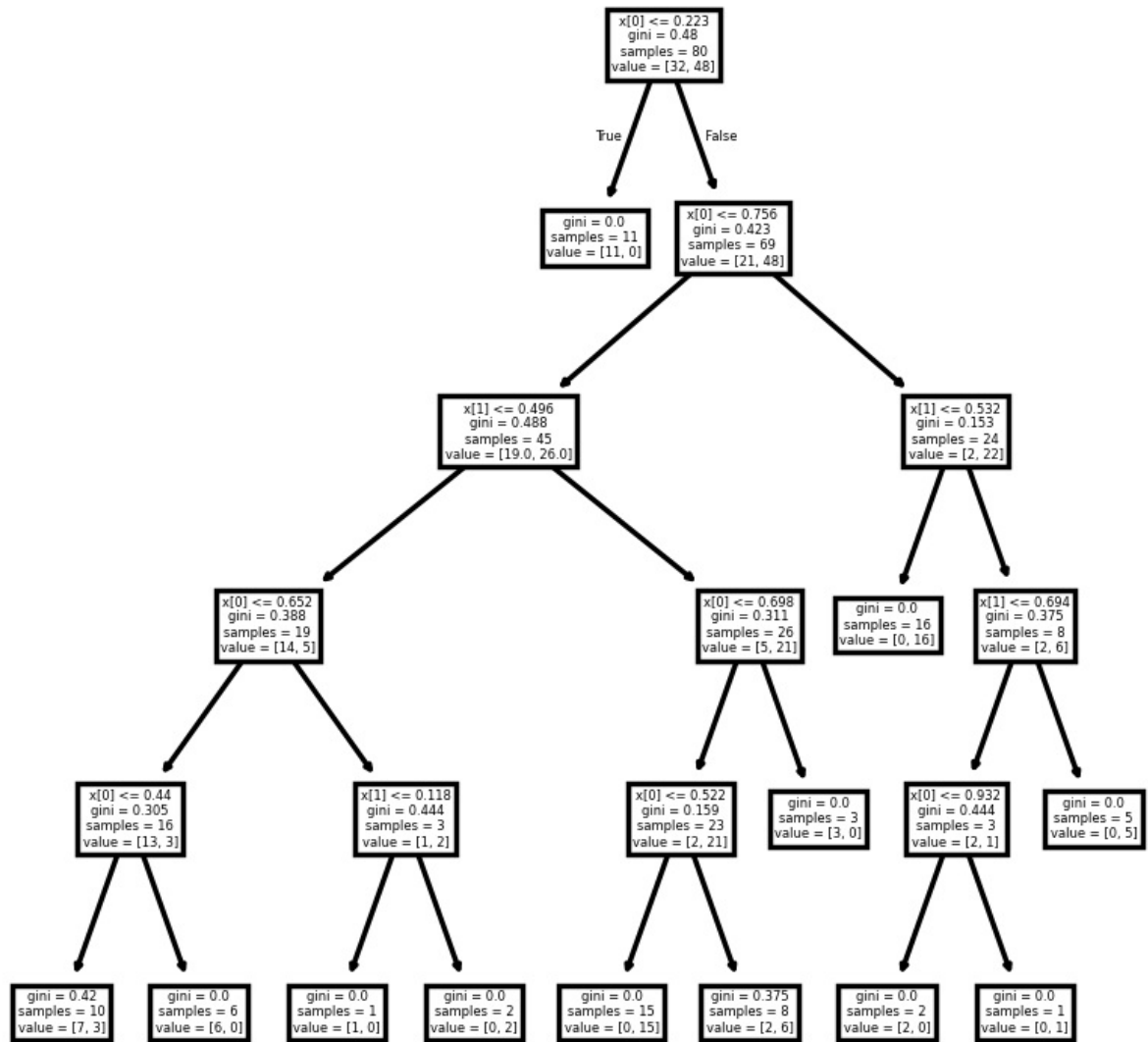
Visualizing the decision tree

The `plot_tree()` function (https://scikit-learn.org/stable/modules/generated/sklearn.tree.plot_tree.html) can generate a simple visualization of your decision tree model. Try out this function below:

```
In [4]: plt.figure(figsize=(4,4),dpi=250)

# YOUR CODE GOES HERE
plot_tree(model)

plt.show()
```



Processing math: 100%

M5-L1 Problem 3

Let's revisit the initial speed vs. launch angle data from the logistic regression module. This time, you will train a decision tree classifier to predict whether a projectile launched with a given speed and angle will hit a target.

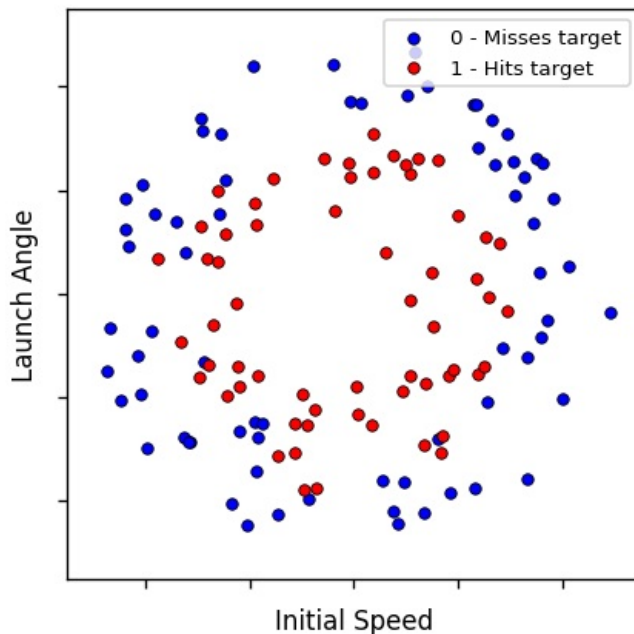
Run this cell to load the data and decision tree tools:

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeClassifier, plot_tree
from matplotlib.colors import ListedColormap

y = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
x1 = np.array([0.02693745, 0.41186575, 0.10363585, 0.08489663, 0.09512868, 0.31121109, 0.16015486, 0.75698706, 0.10363585, 0.08489663, 0.09512868, 0.31121109, 0.16015486, 0.75698706, 0.10363585, 0.08489663, 0.09512868, 0.31121109, 0.16015486, 0.75698706, 0.10363585, 0.08489663, 0.09512868, 0.31121109, 0.16015486, 0.75698706, 0.10363585, 0.08489663, 0.09512868, 0.31121109, 0.16015486, 0.75698706])
x2 = np.array([0.3501823 , 0.10349458, 0.20137442, 0.37973165, 0.71062143, 0.25377085, 0.64055034, 0.29218012, 0.3501823 , 0.10349458, 0.20137442, 0.37973165, 0.71062143, 0.25377085, 0.64055034, 0.29218012, 0.3501823 , 0.10349458, 0.20137442, 0.37973165, 0.71062143, 0.25377085, 0.64055034, 0.29218012, 0.3501823 , 0.10349458, 0.20137442, 0.37973165, 0.71062143, 0.25377085, 0.64055034, 0.29218012])
X = np.vstack([x1, x2]).T

def plot_data(X,y):
    colors=["blue","red"]
    labels = ["0 - Misses target", "1 - Hits target"]
    for i in range(2):
        plt.scatter(X[y==i,0],X[y==i,1],s=20,c=colors[i],edgecolors="black",linewidths=.5,label=labels[i])
        plt.xlabel("Initial Speed")
        plt.ylabel("Launch Angle")
        plt.legend(loc="upper right",prop={'size':8})
    ax = plt.gca()
    ax.set_xticklabels([])
    ax.set_yticklabels([])
    plt.xlim([-0.55,.55])
    plt.ylim([-0.55,.55])

plt.figure(figsize=(4,4),dpi=120)
plot_data(X,y)
plt.show()
```

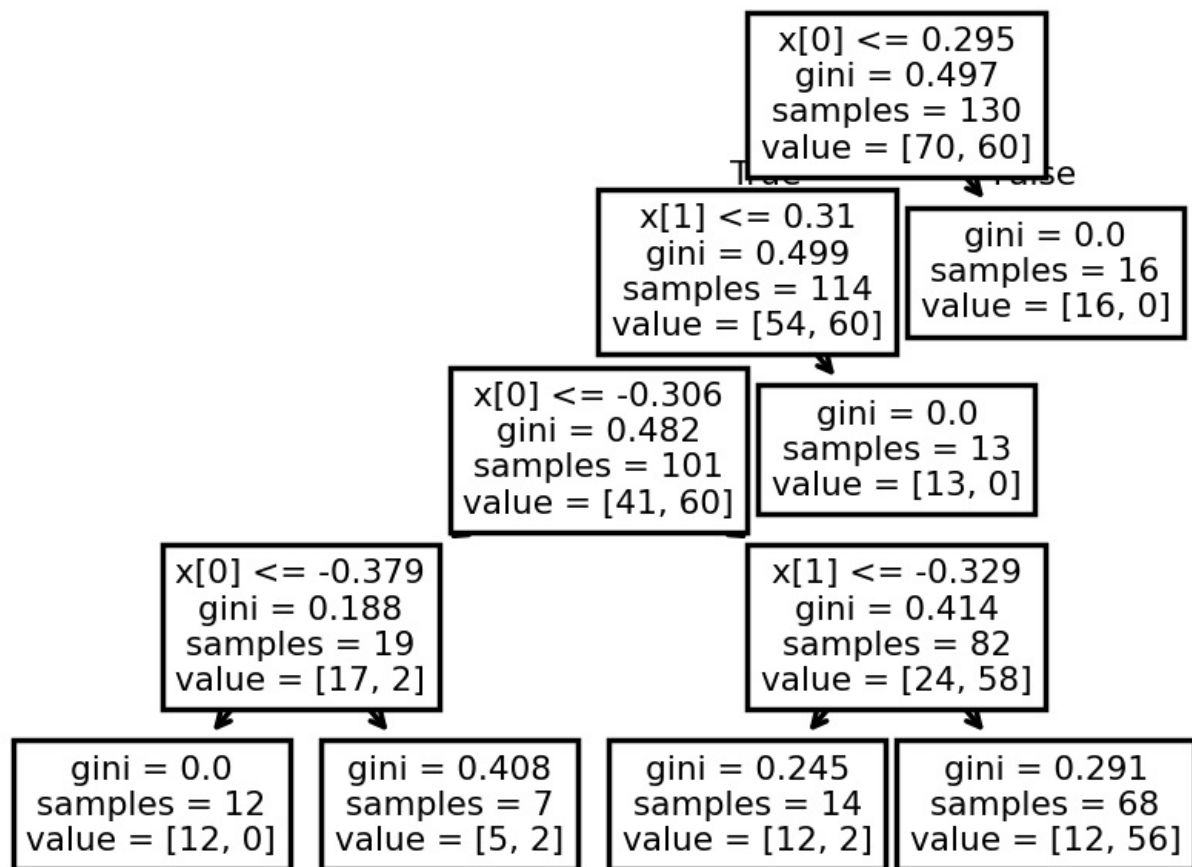


Training a decision tree classifier.

Below, a decision tree of max depth 4 is trained, and the tree is visualized with `plot_tree()`.

```
In [2]: dt = DecisionTreeClassifier(max_depth=4)
dt.fit(X,y)

plt.figure(figsize=(4,3),dpi=250)
plot_tree(dt)
plt.show()
```



Accuracy on training data

Compute the accuracy on the training data with the provided function `get_dt_accuracy(dt, X, y)`. Print the result.

```
In [3]: def get_dt_accuracy(dt, X, y):
        pred = dt.predict(X)
        return 100*np.sum(pred == y)/len(y)

# YOUR CODE GOES HERE
accuracy = get_dt_accuracy(dt, X, y)
print("Accuracy of prediction by decision tree: ", accuracy)
```

Accuracy of prediction by decision tree: 87.6923076923077

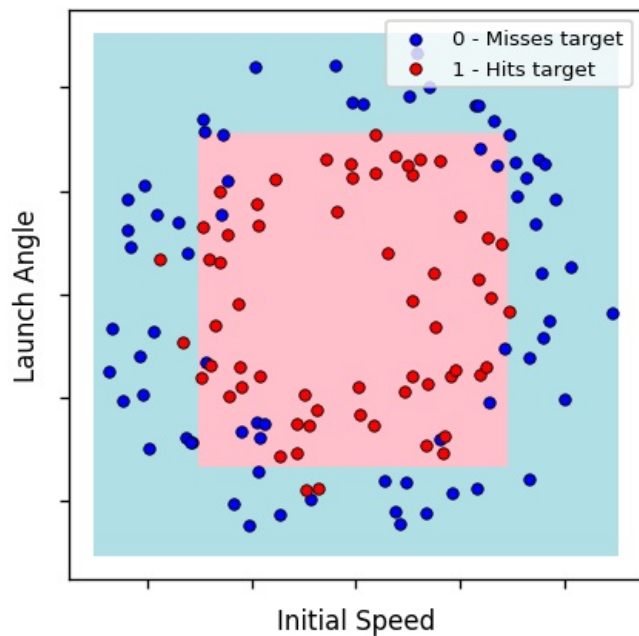
Visualizing tree predictions

By evaluating the model on a meshgrid of results, we can look at how our model performs on the input space:

```
In [4]: vals = np.linspace(-.5,.5,100)
        x1grid, x2grid = np.meshgrid(vals, vals)
        X_test = np.vstack([x1grid.flatten(), x2grid.flatten()]).T

        pred = dt.predict(X_test)

        plt.figure(figsize=(4,4),dpi=120)
        bgcolors = ListedColormap(["powderblue","pink"])
        plt.pcolormesh(x1grid, x2grid, pred.reshape(x1grid.shape), shading="nearest", cmap=bgcolors)
        plot_data(X,y)
        plt.show()
```



Expanded feature set

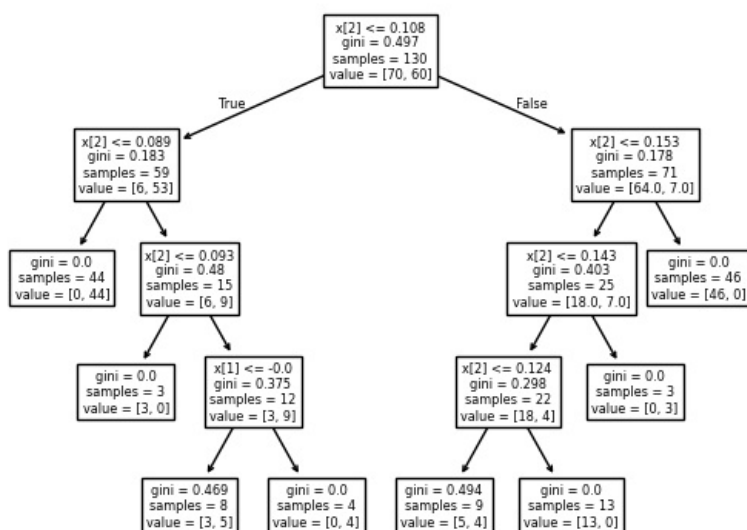
Now, we will add a third feature that (for this problem) happens to be very useful. That feature is $x_1^2 + x_2^2$. A new training input `X_ex` is generated below containing this additional feature.

Train a new decision tree, max depth 4, on this data. Then visualize the tree with `plot_tree()`.

```
In [10]: def feature_expand(X):
x1 = X[:,0].reshape(-1, 1)
x2 = X[:,1].reshape(-1, 1)
columns = [x1, x2, x1*x1 + x2*x2]
return np.concatenate(columns, axis=1)
```

```
X_ex = feature_expand(X)
```

```
# YOUR CODE GOES HERE
# Train a new decision tree on X_ex, y
dt.fit(X_ex, y)
# Plot the tree
plot_tree(dt)
plt.show()
```



Accuracy on training data: expanded features

Compute the accuracy of this new model its training data. It should have increased. Note that the useful features to expand will vary significantly from problem to problem.

```
In [7]: # YOUR CODE GOES HERE
accuracy = get_dt_accuracy(dt, X_ex, y)
```



```
print("Accuracy after training on expanded features:", accuracy)
```

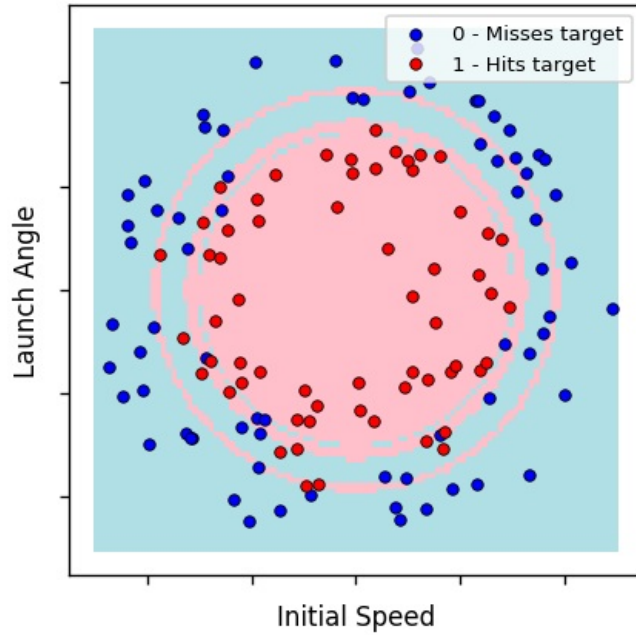
Accuracy after training on expanded features: 94.61538461538461

Visualizing expanded feature results

Use your model to make a prediction called `pred` on the data `X_test_ex`, an expanded meshgrid of points, as indicated. This code will plot the class decisions. Note the difference between this and the previous model, which only had speed and angle as features.

```
In [13]: X_test_ex = feature_expand(X_test)
# YOUR CODE GOES HERE
# Have your model make a prediction, `pred` on X_test_ex
pred = dt.predict(X_test_ex)

plt.figure(figsize=(4,4),dpi=120)
bgcolors = ListedColormap(["powderblue","pink"])
plt.pcolormesh(xlgrid, x2grid, pred.reshape(xlgrid.shape), shading="nearest",cmap=bgcolors)
plot_data(X,y)
plt.show()
```



Processing math: 100%

M5-L2 Problem 1

256 particles of liquid argon are simulated at 100K. A radial distribution function $g(r)$ describes the density of particles a distance of r from each particle in the system. When an $g(r)$ is computed in a simulation, it is done by creating a histogram of particle distances for a single simulation frame, resulting in a noisy function that is most often averaged over several frames.

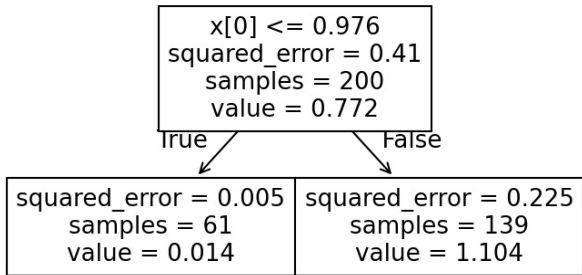
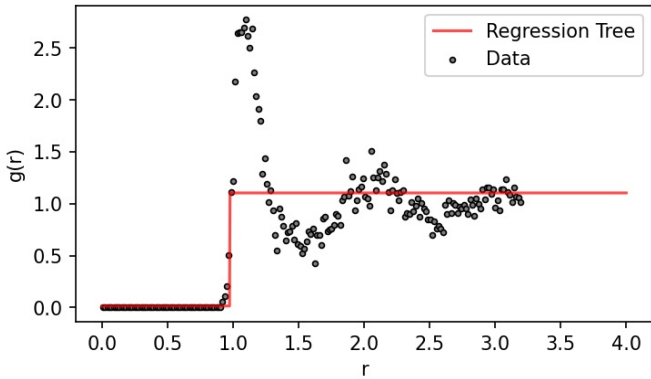
Given $g(r)$ vs. r data for a single frame, you will train a decision tree regressor to represent the underlying function.

First, run the cell below to load the data, etc.:

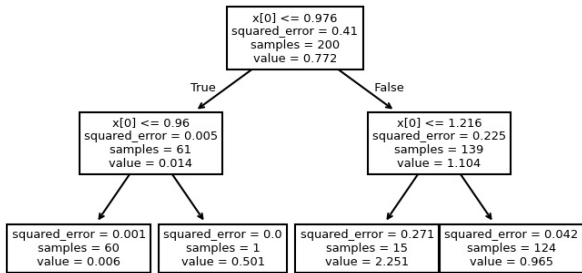
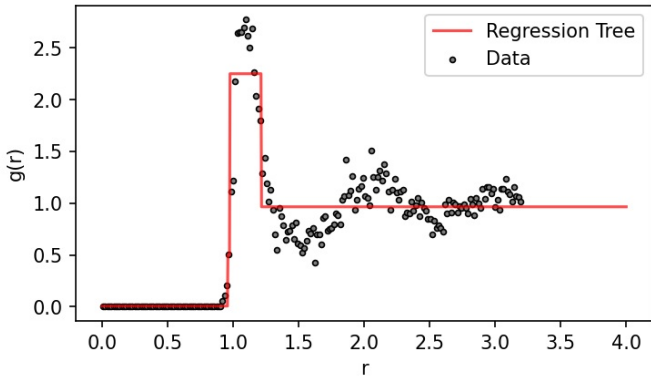
```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeRegressor, plot_tree

r = np.array([0.008,0.024,0.04,0.056,0.072,0.088,0.104,0.12,0.136,0.152,0.168,0.184,0.2,0.216,0.232,0.248,0.264,
0.28,0.296,0.312,0.328,0.344,0.36,0.376,0.392,0.408,0.424,0.44,0.456,0.472,0.488,0.504,0.52,0.536,0.552,0.568,0.584,0.6,0.616,0.632,0.648,0.664,0.68,0.696,0.712,0.728,0.744,0.76,0.776,0.792,0.808,0.824,0.84,0.856,0.872,0.888,0.904,0.92,0.936,0.952,0.968,0.984,1.0,1.016,1.032,1.048,1.064,1.08,1.096,1.112,1.128,1.144,1.16,1.176,1.192,1.208,1.224,1.24,1.256,1.272,1.288,1.304,1.32,1.336,1.352,1.368,1.384,1.4,1.416,1.432,1.448,1.464,1.48,1.496,1.512,1.528,1.544,1.56,1.576,1.592,1.608,1.624,1.64,1.656,1.672,1.688,1.704,1.72,1.736,1.752,1.768,1.784,1.8,1.816,1.832,1.848,1.864,1.88,1.896,1.912,1.928,1.944,1.96,1.976,1.992,2.008,2.024,2.04,2.056,2.072,2.088,2.104,2.12,2.136,2.152,2.168,2.184,2.2,2.216,2.232,2.248,2.264,2.28,2.296,2.312,2.328,2.344,2.36,2.376,2.392,2.408,2.424,2.44,2.456,2.472,2.488,2.504,2.52,2.536,2.552,2.568,2.584,2.6,2.616,2.632,2.648,2.664,2.68,2.696,2.712,2.728,2.744,2.76,2.776,2.792,2.808,2.824,2.84,2.856,2.872,2.888,2.904,2.92,2.936,2.952,2.968,2.984,3.0,3.016,3.032,3.048,3.064,3.08,3.096,3.112,3.128,3.144,3.16,3.176,3.192,3.208,3.224,3.24,3.256,3.272,3.288,3.304,3.32,3.336,3.352,3.368,3.384,3.4,3.416,3.432,3.448,3.464,3.48,3.496,3.512,3.528,3.544,3.56,3.576,3.592,3.608,3.624,3.64,3.656,3.672,3.688,3.704,3.72,3.736,3.752,3.768,3.784,3.8,3.816,3.832,3.848,3.864,3.88,3.896,3.912,3.928,3.944,3.96,3.976,3.992,4.008,4.024,4.04,4.056,4.072,4.088,4.104,4.12,4.136,4.152,4.168,4.184,4.2,4.216,4.232,4.248,4.264,4.28,4.296,4.312,4.328,4.344,4.36,4.376,4.392,4.408,4.424,4.44,4.456,4.472,4.488,4.504,4.52,4.536,4.552,4.568,4.584,4.6,4.616,4.632,4.648,4.664,4.68,4.696,4.712,4.728,4.744,4.76,4.776,4.792,4.808,4.824,4.84,4.856,4.872,4.888,4.904,4.92,4.936,4.952,4.968,4.984,5.0,5.016,5.032,5.048,5.064,5.08,5.096,5.112,5.128,5.144,5.16,5.176,5.192,5.208,5.224,5.24,5.256,5.272,5.288,5.304,5.32,5.336,5.352,5.368,5.384,5.4,5.416,5.432,5.448,5.464,5.48,5.496,5.512,5.528,5.544,5.56,5.576,5.592,5.608,5.624,5.64,5.656,5.672,5.688,5.704,5.72,5.736,5.752,5.768,5.784,5.8,5.816,5.832,5.848,5.864,5.88,5.896,5.912,5.928,5.944,5.96,5.976,5.992,6.008,6.024,6.04,6.056,6.072,6.088,6.104,6.12,6.136,6.152,6.168,6.184,6.2,6.216,6.232,6.248,6.264,6.28,6.296,6.312,6.328,6.344,6.36,6.376,6.392,6.408,6.424,6.44,6.456,6.472,6.488,6.504,6.52,6.536,6.552,6.568,6.584,6.6,6.616,6.632,6.648,6.664,6.68,6.696,6.712,6.728,6.744,6.76,6.776,6.792,6.808,6.824,6.84,6.856,6.872,6.888,6.904,6.92,6.936,6.952,6.968,6.984,7.0,7.016,7.032,7.048,7.064,7.08,7.096,7.112,7.128,7.144,7.16,7.176,7.192,7.208,7.224,7.24,7.256,7.272,7.288,7.304,7.32,7.336,7.352,7.368,7.384,7.4,7.416,7.432,7.448,7.464,7.48,7.496,7.512,7.528,7.544,7.56,7.576,7.592,7.608,7.624,7.64,7.656,7.672,7.688,7.704,7.72,7.736,7.752,7.768,7.784,7.8,7.816,7.832,7.848,7.864,7.88,7.896,7.912,7.928,7.944,7.96,7.976,7.992,8.008,8.024,8.04,8.056,8.072,8.088,8.104,8.12,8.136,8.152,8.168,8.184,8.2,8.216,8.232,8.248,8.264,8.28,8.296,8.312,8.328,8.344,8.36,8.376,8.392,8.408,8.424,8.44,8.456,8.472,8.488,8.504,8.52,8.536,8.552,8.568,8.584,8.6,8.616,8.632,8.648,8.664,8.68,8.696,8.712,8.728,8.744,8.76,8.776,8.792,8.808,8.824,8.84,8.856,8.872,8.888,8.904,8.92,8.936,8.952,8.968,8.984,9.0,9.016,9.032,9.048,9.064,9.08,9.096,9.112,9.128,9.144,9.16,9.176,9.192,9.208,9.224,9.24,9.256,9.272,9.288,9.304,9.32,9.336,9.352,9.368,9.384,9.4,9.416,9.432,9.448,9.464,9.48,9.496,9.512,9.528,9.544,9.56,9.576,9.592,9.608,9.624,9.64,9.656,9.672,9.688,9.704,9.72,9.736,9.752,9.768,9.784,9.8,9.816,9.832,9.848,9.864,9.88,9.896,9.912,9.928,9.944,9.96,9.976,9.992,10.008,10.024,10.04,10.056,10.072,10.088,10.104,10.12,10.136,10.152,10.168,10.184,10.2,10.216,10.232,10.248,10.264,10.28,10.296,10.312,10.328,10.344,10.36,10.376,10.392,10.408,10.424,10.44,10.456,10.472,10.488,10.504,10.52,10.536,10.552,10.568,10.584,10.6,10.616,10.632,10.648,10.664,10.68,10.696,10.712,10.728,10.744,10.76,10.776,10.792,10.808,10.824,10.84,10.856,10.872,10.888,10.904,10.92,10.936,10.952,10.968,10.984,11.0,11.016,11.032,11.048,11.064,11.08,11.096,11.112,11.128,11.144,11.16,11.176,11.192,11.208,11.224,11.24,11.256,11.272,11.288,11.304,11.32,11.336,11.352,11.368,11.384,11.4,11.416,11.432,11.448,11.464,11.48,11.496,11.512,11.528,11.544,11.56,11.576,11.592,11.608,11.624,11.64,11.656,11.672,11.688,11.704,11.72,11.736,11.752,11.768,11.784,11.8,11.816,11.832,11.848,11.864,11.88,11.896,11.912,11.928,11.944,11.96,11.976,11.992,12.008,12.024,12.04,12.056,12.072,12.088,12.104,12.12,12.136,12.152,12.168,12.184,12.2,12.216,12.232,12.248,12.264,12.28,12.296,12.312,12.328,12.344,12.36,12.376,12.392,12.408,12.424,12.44,12.456,12.472,12.488,12.504,12.52,12.536,12.552,12.568,12.584,12.6,12.616,12.632,12.648,12.664,12.68,12.696,12.712,12.728,12.744,12.76,12.776,12.792,12.808,12.824,12.84,12.856,12.872,12.888,12.904,12.92,12.936,12.952,12.968,12.984,13.0,13.016,13.032,13.048,13.064,13.08,13.096,13.112,13.128,13.144,13.16,13.176,13.192,13.208,13.224,13.24,13.256,13.272,13.288,13.304,13.32,13.336,13.352,13.368,13.384,13.4,13.416,13.432,13.448,13.464,13.48,13.496,13.512,13.528,13.544,13.56,13.576,13.592,13.608,13.624,13.64,13.656,13.672,13.688,13.704,13.72,13.736,13.752,13.768,13.784,13.8,13.816,13.832,13.848,13.864,13.88,13.896,13.912,13.928,13.944,13.96,13.976,13.992,14.008,14.024,14.04,14.056,14.072,14.088,14.104,14.12,14.136,14.152,14.168,14.184,14.2,14.216,14.232,14.248,14.264,14.28,14.296,14.312,14.328,14.344,14.36,14.376,14.392,14.408,14.424,14.44,14.456,14.472,14.488,14.504,14.52,14.536,14.552,14.568,14.584,14.6,14.616,14.632,14.648,14.664,14.68,14.696,14.712,14.728,14.744,14.76,14.776,14.792,14.808,14.824,14.84,14.856,14.872,14.888,14.904,14.92,14.936,14.952,14.968,14.984,15.0,15.016,15.032,15.048,15.064,15.08,15.096,15.112,15.128,15.144,15.16,15.176,15.192,15.208,15.224,15.24,15.256,15.272,15.288,15.304,15.32,15.336,15.352,15.368,15.384,15.4,15.416,15.432,15.448,15.464,15.48,15.496,15.512,15.528,15.544,15.56,15.576,15.592,15.608,15.624,15.64,15.656,15.672,15.688,15.704,15.72,15.736,15.752,15.768,15.784,15.8,15.816,15.832,15.848,15.864,15.88,15.896,15.912,15.928,15.944,15.96,15.976,15.992,16.008,16.024,16.04,16.056,16.072,16.088,16.104,16.12,16.136,16.152,16.168,16.184,16.2,16.216,16.232,16.248,16.264,16.28,16.296,16.312,16.328,16.344,16.36,16.376,16.392,16.408,16.424,16.44,16.456,16.472,16.488,16.504,16.52,16.536,16.552,16.568,16.584,16.6,16.616,16.632,16.648,16.664,16.68,16.696,16.712,16.728,16.744,16.76,16.776,16.792,16.808,16.824,16.84,16.856,16.872,16.888,16.904,16.92,16.936,16.952,16.968,16.984,17.0,17.016,17.032,17.048,17.064,17.08,17.096,17.112,17.128,17.144,17.16,17.176,17.192,17.208,17.224,17.24,17.256,17.272,17.288,17.304,17.32,17.336,17.352,17.368,17.384,17.4,17.416,17.432,17.448,17.464,17.48,17.496,17.512,17.528,17.544,17.56,17.576,17.592,17.608,17.624,17.64,17.656,17.672,17.688,17.704,17.72,17.736,17.752,17.768,17.784,17.8,17.816,17.832,17.848,17.864,17.88,17.896,17.912,17.928,17.944,17.96,17.976,17.992,18.008,18.024,18.04,18.056,18.072,18.088,18.104,18.12,18.136,18.152,18.168,18.184,18.2,18.216,18.232,18.248,18.264,18.28,18.296,18.312,18.328,18.344,18.36,18.376,18.392,18.408,18.424,18.44,18.456,18.472,18.488,18.504,18.52,18.536,18.552,18.568,18.584,18.6,18.616,18.632,18.648,18.664,18.68,18.696,18.712,18.728,18.744,18.76,18.776,18.792,18.808,18.824,18.84,18.856,18.872,18.888,18.904,18.92,18.936,18.952,18.968,18.984,19.0,19.016,19.032,19.048,19.064,19.08,19.096,19.112,19.128,19.144,19.16,19.176,19.192,19.208,19.224,19.24,19.256,19.272,19.288,19.304,19.32,19.336,19.352,19.368,19.384,19.4,19.416,19.432,19.448,19.464,19.48,19.496,19.512,19.528,19.544,19.56,19.576,19.592,19.608,19.624,19.64,19.656,19.672,19.688,19.704,19.72,19.736,19.752,19.768,19.784,19.8,19.816,19.832,19.848,19.864,19.88,19.896,19.912,19.928,19.944,19.96,19.976,19.992,20.008,20.024,20.04,20.056,20.072,20.088,20.104,20.12,20.136,20.152,20.168,20.184,20.2,20.216,20.232,20.248,20.264,20.28,20.296,20.312,20.328,20.344,20.36,20.376,20.392,20.408,20.424,20.44,20.456,20.472,20.488,20.504,20.52,20.536,20.552,20.568,20.584,20.6,20.616,20.632,20.648,20.664,20.68,20.696,20.712,20.728,20.744,20.76,20.776,20.792,20.808,20.824,20.84,20.856,20.872,20.888,20.904,20.92,20.936,20.952,20.968,20.984,21.0,21.016,21.032,21.048,21.064,21.08,21.096,21.112,21.128,21.144,21.16,21.176,21.192,21.208,21.224,21.24,21.256,21.272,21.288,21.304,21.32,21.336,21.352,21.368,21.384,21.4,21.416,21.432,21.448,21.464,21.48,21.496,21.512,21.528,21.544,21.56,21.576,21.592,21.608,21.624,21.64,21.656,21.672,21.688,21.704,21.72,21.736,21.752,21.768,21.784,21.8,21.816,21.832,21.848,21.864,21.88,21.896,21.912,21.928,21.944,21.96,21.976,21.992,22.008,22.024,22.04,22.056,22.072,22.088,22.104,22.12,22.136,22.152,22.168,22.184,22.2,22.216,22.232,22.248,22.264,22.28,22.296,22.312,22.328,22.344,22.36,22.376,22.392,22.408,22.424,22.44,22.456,22.472,22.488,22.504,22.52,22.536,22.552,22.568,22.584,22.6,22.616,22.632,22.648,22.664,22.68,22.696,22.712,22.728,22.744,22.76,22.776,22.792,22.808,22.824,22.84,22.856,22.872,22.888,22.904,22.92,22.936,22.952,22.968,22.984,23.0,23.016,23.032,23.048,23.064,23.08,23.096,23.112,23.128,23.144,23.16,23.176,23.192,23.208,23.224,23.24,23.256,23.272,23.288,23.304,23.32,23.336,23.352,23.368,23.384,23.4,23.416,23.432,23.448,23.464,23.48,23.496,23.512,23.528,23.544,23.56,23.576,23.592,23.608,23.624,23.64,23.656,23.672,23.688,23.704,23.72,23.736,23.752,23.768,23.784,23.8,23.816,23.832,23.848,23.864,23.88,23.896,23.912,23.928,23.944,23.96,23.976,23.992,24.008,24.024,24.04,24.056,24.072,24.088,24.104,24.12,24.136,24.152,24.168,24.184,24.2,24.216,24.232,24.248,24.264,24.28,24.296,24.312,24.328,24.344,24.36,24.376,24.392,24.408,24.424,24.44,24.456,24.472,24.488,24.504,24.52,24.536,24.552,24.568,24.584,24.6,24.616,24.632,24.648,24.664,24.68,24.696,24.712,24.728,24.744,24.76,24.776,24.792,24.808,24.824,24.84,24.856,24.872,24.888,24.904,24.92,24.936,24.952,24.968,24.984,25.0,25.016,25.032,25.048,25.064,25.08,25.096,25.112,25.128,25.144,25.16,25.176,25.192,25.208,25.224,25.24,25.256,25.272,25.288,25.304,25.32,25.336,25.352,25.368,25.384,25.4,25.416,25.432,25.448,25.464,25.48,25.496,25.512,25.528,25.544,25.56,25.576,25.592,25.608,25.624,25.64,25.656,25.672,25.688,25.704,25.72,25.736,25.752,25.768,25.784,25.8,25.816,25.832,25.848,25.864,25.88,25.896,25.912,25.928,25.944,25.96,25.976,25.992,26.008,26.024,26.04,26.056,26.072,26.088,26.104,26.12,26.136,26.152,26.168,26.184,26.2,26.216,26.232,26.248,26.264,26.28,26.296,26.312,26.328,26.344,26.36,26.376,26.392,26.408,26.424,26.44,26.456,26.472,26.488,26.504,26.52,26.536,26.552,26.568,26.584,26.6,26.616,26.632,26.648,26.664,26.68,26.696,26.712,26.728,26.744,26.76,26.776,26.792,26.808,26.824,26.84,26.856,26.872,26.888,26.904,26.92,26.936,26.952,26.968,26.984,27.0,27.016,27.032,27.048,27.064,27.08,27.096,27.112,27.128,27.144,27.16,27.176,27.192,27.208,27.224,27.24,27.256,27.272,27.288,27.304,27.32,27.336,27.352,27.368,27.384,27.4,27.416,27.432,27.448,27.464,27.48,27.496,27.512,27.528,27.544,27.56,27.576,27.592,27.608,27.624,27.64,27.656,27.672,27.688,27.704,27.72,27.736,27.752,27.768,27.784,27.8,27.816,27.832,27.848,27.864,27.88,27.896,27.912,27.928,27.944,27.96,27.976,27.992,28.008,28.024,28.04,28.056,28.072,28.088,28.104,28.12,28.136,28.152,28.168,28.184,28.2,28.216,28.232,28.248,28.264,28.28,28.296,28.312,28.328,28.344,28.36,28.376,28.392,28.408,28.424,28.44,28.456,28.472,28.488,28.504,28.52,28.536,28.552,28.568,28.584,28.6,28.616,28.632,28.648,28.664,28.68,28.696,28.712,28.728,28.744,28.76,28.776,28.792,28.808,28.824,28.84,28.856,28.872,28.888,28.904,28.92,28.936,28.952,28.968,28.984,29.0,29.016,29.032,29.048,29.064,29.08,29.096,29.112,29.128,29.144,29.16,29.176,29.192,29.208,29.224,29.24,29.256,29.272,29.288,29.304,29.32,29.336,29.352,29.368,29.384,29.4,29.416,29.432,29.448,29.464,29.48,29.496,29.512,29.528,29.544,29.56,29.576,29.592,29.608,29.624,29.64,29.656,29.672,29.688,29.704,29.72,29.736,29.752,29.768,29.784,29.8,29.816,29.832,29.848,29.864,29.88,29.896,29.912,29.928,29.944,29.96,29.976,
```

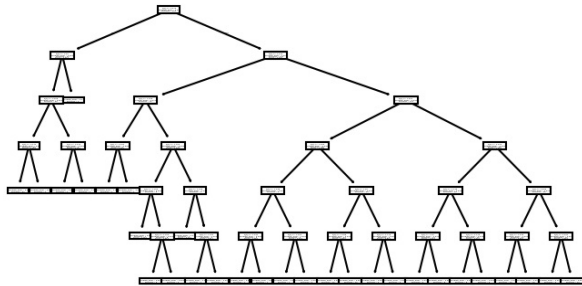
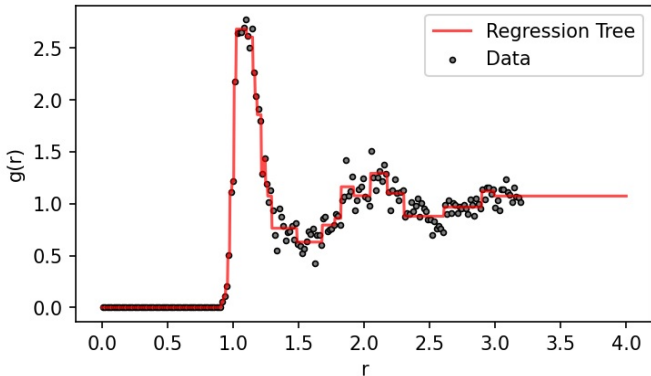
```
for max_depth in [1, 2, 6, 10]:
    # YOUR CODE GOES HERE
    # Create and fit `dt`
    dt = DecisionTreeRegressor(max_depth = max_depth)
    r_ = r.reshape(-1, 1)
    dt.fit(r_, g)
    plot(r,g,dt)
```



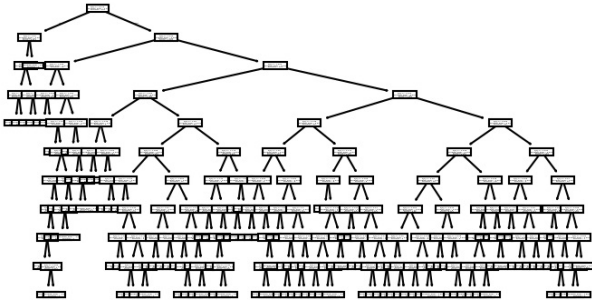
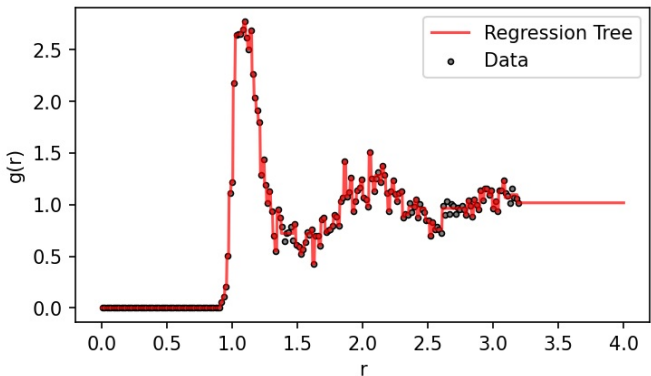
Tree max. depth: 1



Tree max. depth: 2



Tree max. depth: 6



Tree max. depth: 10

Processing math: 100%

M5-L2 Problem 2

Stress-strain measurements have been collected for many samples across many parts, resulting in much noisier data than would come from a tensile test, for example. Your job is to train an ensemble of decision trees that can predict stress for an input strain.

Scikit-Learn's `RandomForestRegressor()` has several parameters that you will experiment with below.

Run each cell; then, experiment with different settings of the `RandomForestRegressor()` to answer the questions at the end.

```
In [1]: # Import libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.ensemble import RandomForestRegressor
%matplotlib inline
from ipywidgets import interact, interactive, fixed, interact_manual, Layout, FloatSlider, Dropdown

# Load the data
y = np.array([133.18473289, 366.12422297, 453.70990214, 479.37136253, 238.16361712, 39.91719443, 282.21638562,
x = np.array([0.47358185, 0.80005535, 1.10968143, 1.85282726, 0.58177792, 0.24407275, 0.67817621, 0.59768343, 1

In [2]: def plot(n_estimators, max_leaf_nodes, bootstrap):
    n_estimators = [1,10,20,30,40,50,60,70,80,90,100][int(n_estimators)]
    max_leaf_nodes = int(max_leaf_nodes)
    model = RandomForestRegressor(n_estimators=n_estimators,
                                bootstrap=(True if "On" in bootstrap else False),
                                max_leaf_nodes=max_leaf_nodes,
                                random_state=0)

    model.fit(x.reshape(-1,1), y)

    xs = np.linspace(min(x),max(x),500)
    ys = model.predict(xs.reshape(-1,1))

    plt.figure(figsize=(5,3),dpi=150)
    plt.scatter(x,y,s=20,color="cornflowerblue",edgecolor="navy",label="Data")
    plt.plot(xs, ys, c="red",linewidth=2,label="Mean prediction")
    for i,dt in enumerate(model.estimators_):
        label = "Tree predictions" if i == 0 else None
        plt.plot(xs, dt.predict(xs.reshape(-1,1)), c="gray",linewidth=.5,zorder=-1, label = label)

    plt.legend(loc="lower right",prop={"size":8})
    plt.xlabel("Strain, %")
    plt.ylabel("Stress, MPa")
    plt.title(f"Num. estimators: {n_estimators}, Max leaves = {max_leaf_nodes}, Bootstrapping: {bootstrap}",font
    plt.show()

slider1 = FloatSlider(
    value=2,
    min=0,
    max=10,
    step=1,
    description='# Estimators',
    disabled=False,
    continuous_update=True,
    orientation='horizontal',
    readout=False,
    layout = Layout(width='550px')
)

slider2 = FloatSlider(
    value=5,
    min=2,
    max=25,
    step=1,
    description='Max Leaves',
    disabled=False,
    continuous_update=True,
    orientation='horizontal',
    readout=False,
    layout = Layout(width='550px')
)

dropdown = Dropdown(
    options=["On (66% of data)", "Off"],
    value="On (66% of data)",
    description='Bootstrap',
    disabled=False,
)
```

```

interactive_plot = interactive(
    plot,
    bootstrap = dropdown,
    n_estimators = slider1,
    max_leaf_nodes = slider2
)
output = interactive_plot.children[-1]
output.layout.height = '500px'

interactive_plot

```

Out[2]: interactive(children=(FloatSlider(value=2.0, description='# Estimators', layout=Layout(width='550px'), max=10...

Questions

1. Keep bootstrapping on and set max leaf nodes constant at 3. Describe what happens to the mean prediction as the number of estimators increases.
 2. Keep bootstrapping on and set number of estimators constant at 100. Describe what happens to the mean prediction as the leaf node maximum increases.
 3. Now disable bootstrapping. Notice that all of the predictions are the same -- the gray lines are behind the red. Why is this? (Hint: Think about the number of features in this dataset.)
-
1. As the number of estimators increases while the number of max leaf nodes remain constant at 3, the mean prediction becomes more stable as the fluctuations in the individual trees average out. However, since each tree is constrained to 3 nodes, their complexity is still limited and they maintain a high bias. The mean prediction ends up being rather smooth and misses out on capturing finer details.
 2. Increasing the maximum number of nodes in individual trees increases their complexities, and this helps in capturing more complex data. The mean prediction will have reduced bias, and maintain a somewhat low variance due to averaging. As the number of nodes increases too much, individual tree may end up overfitting to their samples. This will cause it to be susceptible to outliers and noise. The mean prediction still somewhat converges thanks to averaging, however, the change in performance if any, might not be significant.
 3. Since the dataset has very few features, each tree has only limited options for splitting the data between its nodes. So it's highly likely that the algorithm makes the same split across all the trees. This means that the trees are essentially identical clones of each other. So, since all the predictions are the same, they appear to be hidden behind the mean prediction since they overlap perfectly.

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

Problem 6

Problem Description

In this problem you will train decision tree and random forest models using sklearn on a real world dataset. The dataset is the *Cylinder Bands Data Set* from the UCI Machine Learning Repository: <https://archive.ics.uci.edu/ml/datasets/Cylinder+Bands>. The dataset is generated from rotogravure printers, with 39 unique features, and a binary classification label for each sample. The class is either 0, for 'band' or 1 for 'no band', where banding is an undesirable process delay that arises during the rotogravure printing process. By training ML models on this dataset, you could help identify or predict cases where these process delays are avoidable, thereby improving the efficiency of the printing. For the sake of this exercise, we only consider features 21-39 in the above link, and have removed any samples with missing values in that range. No further processing of the data is required on your behalf. The data has been partitioned into a training and testing set using an 80/20 split. Your models will be trained on just the train set, and accuracy results will be reported on both the training and testing sets.

Fill out the notebook as instructed, making the requested plots and printing necessary values.

You are welcome to use any of the code provided in the lecture activities.

Summary of deliverables:

- Accuracy function
- Report accuracy of the DT model on the training and testing set
- Report accuracy of the Random Forest model on the training and testing set

Imports and Utility Functions:

```
In [1]: import numpy as np
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
```

Load the data

Use the `np.load()` function to load "w5-hw1-train.npy" (training data) and "w5-hw1-test.npy" (testing data). The first 19 columns of each are the features. The last column is the label

```
In [12]: # YOUR CODE GOES HERE
train_data = np.load("data/w5-hw1-train.npy")
test_data = np.load("data/w5-hw1-test.npy")

X_train = train_data[:, :19]
y_train = train_data[:, -1]

X_test = test_data[:, :19]
y_test = test_data[:, -1]
```

Write an accuracy function

Write a function `accuracy(pred, label)` that takes in the models prediction, and returns the percentage of predictions that match the corresponding labels.

```
In [9]: # YOUR CODE GOES HERE
def accuracy(pred, label):
    acc = np.sum((pred == label)) * 100 / len(pred)
    return acc
```

Train a decision tree model

Train a decision tree using `DecisionTreeClassifier()` with a `max_depth` of 10 and using a `random_state` of 0 to ensure repeatable results. Print the accuracy of the model on both the training and testing sets.

```
In [19]: # YOUR CODE GOES HERE
model = DecisionTreeClassifier(max_depth = 10, random_state = 0)
model.fit(X_train, y_train)

print("Model accuracy of DecisionTreeClassifier on training data: ", accuracy(model.predict(X_train), y_train))
print("Model accuracy of DecisionTreeClassifier on test data: ", accuracy(model.predict(X_test), y_test))
```

Model accuracy on training data: 93.12714776632302
Model accuracy on test data: 65.75342465753425

Train a random forest model

Train a random forest model using `RandomForestClassifier()` with a `max_depth` of 10, a `n_estimators` of 100, and using a random state of `0` to ensure repeatable results. Print the accuracy of the model on both the training and testing sets.

```
In [20]: # YOUR CODE GOES HERE
model_ = RandomForestClassifier(max_depth = 10, n_estimators = 100, random_state = 0)
model_.fit(X_train, y_train)

print("Model accuracy of RandomForestClassifier on training data: ", accuracy(model_.predict(X_train), y_train))
print("Model accuracy of RandomForestClassifier on test data: ", accuracy(model_.predict(X_test), y_test))
```

Model accuracy of RandomForestClassifier on training data: 100.0
Model accuracy of RandomForestClassifier on test data: 82.1917808219178

Discuss the performance of the models

Compare the training and testing accuracy of the two models, and explain why the random forest model is advantageous compared to a standard decision tree model

A single decision tree might not be very accurate, especially if the depth is limited, there is noise in the data etc. Whereas in a random forest model, predictions of multiple decision trees are averaged out, resulting in irregularities being smoothed out and making the model less sensitive to noise and outliers. A decision might also be overfit to the data, which is avoided in random forest models, as the averaging and randomness generalizes the data better.

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

Problem 2

Problem Description

In this problem, you are given a dataset with two input features and one output. You will use a regression tree to make predictions for this data, evaluating each model on both training and testing data. Then, you will repeat this for multiple random forests.

Fill out the notebook as instructed, making the requested plots and printing necessary values.

You are welcome to use any of the code provided in the lecture activities.

Summary of deliverables:

- RMSE function
- Create 4 decision tree prediction surface plots
- Create 4 random forest prediction surface plots
- Print RMSE for train and test data for 4 decision tree models
- Print RMSE for train and test data for 4 random forest models
- Answer the 3 questions posed throughout

Imports and Utility Functions:

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from mpl_toolkits.mplot3d import Axes3D
from matplotlib import cm

def make_plot(X,y,model, title=""):
    res = 100
    xrange = np.linspace(min(X[:,0]),max(X[:,0]),res)
    yrange = np.linspace(min(X[:,1]),max(X[:,1]),res)
    x1,x2 = np.meshgrid(xrange,yrange)
    xmesh = np.vstack([x1.flatten(),x2.flatten()]).T
    z = model.predict(xmesh).reshape(res,res)

    fig = plt.figure(figsize=(12,10))
    plt.subplots_adjust(left=0.3,right=0.9,bottom=.3,top=.9)
    ax = fig.add_subplot(111, projection='3d')
    ax.plot_surface(x1,x2,z,cmap=cm.coolwarm,linewidth=0,alpha=0.9)
    ax.scatter(X[:,0],X[:,1],y,'o',c='black')
    ax.set_xlabel('$x_1$')
    ax.set_ylabel('$x_2$')
    ax.set_zlabel('$y$')
    plt.title(title)
    plt.show()
```

Load the data

Use the `np.load()` function to load "w5-hw2-train.npy" (training data) and "w5-hw2-test.npy" (testing data). The first two columns of each are the input features. The last column is the output. You should end up with 4 variables, input and output for each of the datasets.

```
In [4]: # YOUR CODE GOES HERE
data_train = np.load("data/w5-hw2-train.npy")
data_test = np.load("data/w5-hw2-test.npy")

X_train = data_train[:, 0:-1]
y_train = data_train[:, -1]

X_test = data_test[:, 0:-1]
y_test = data_test[:, -1]
```

RMSE function

Complete a root-mean-squared-error function, `RMSE(y, pred)`, which takes in two arrays, and computes the RMSE between them:

```
In [6]: def RMSE(y, pred):
# YOUR CODE GOES HERE
rms = np.sqrt(np.sum(np.square(pred - y)) / len(y))
return rms
```


Regression trees

Train 4 regression trees in sklearn, with max depth values [2,5,10,25]. Train your models on the training data.

Plot the predictions as a surface plot along with test points -- you can use the provided function: `make_plot(X, y, model, title)`.

For each model, compute the train and test RMSE by calling your RMSE function. Print these results.

```
In [13]: # YOUR CODE GOES HERE
depth = [2, 5, 10, 25]

for i in depth:
    dt = DecisionTreeRegressor(max_depth = i)
    dt.fit(X_train, y_train)

    pred_train = dt.predict(X_train)
    pred_test = dt.predict(X_test)
    rms_train = RMSE(y_train, pred_train)
    rms_test = RMSE(y_test, pred_test)

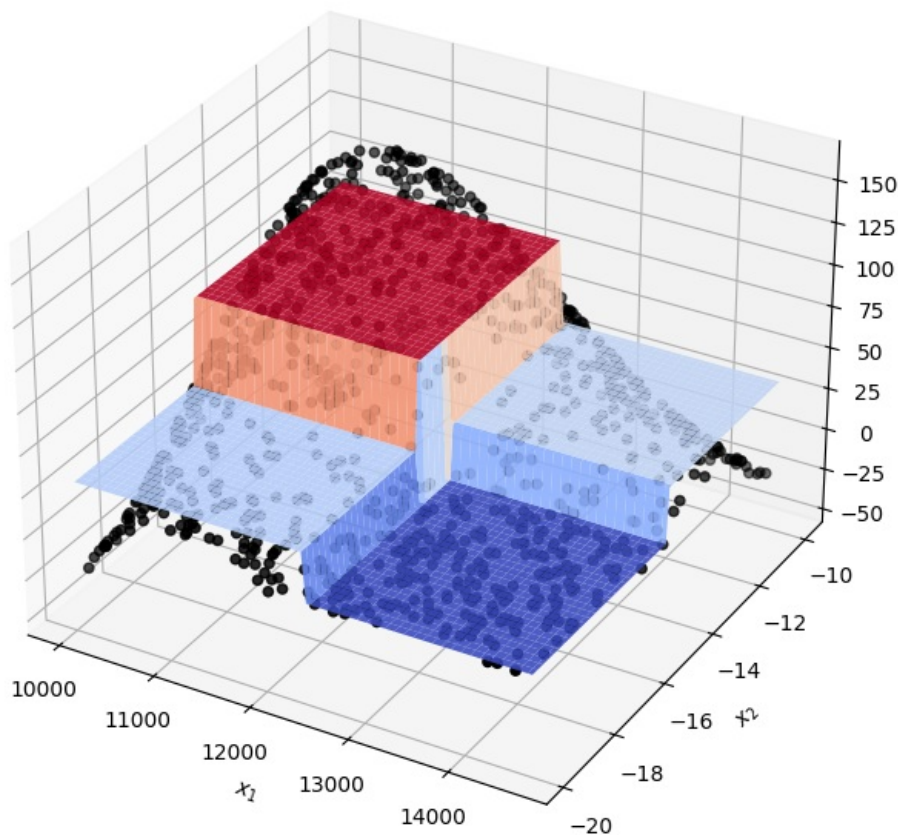
    print("\nMax Depth: ", i)
    print("RMS error for Decision Tree Regressor on training data: ", rms_train)
    print("RMS error for Decision Tree Regressor on test data: ", rms_test)
    make_plot(X_train, y_train, dt, title = "training data")
    make_plot(X_test, y_test, dt, title = "test data")
```

Max Depth: 2

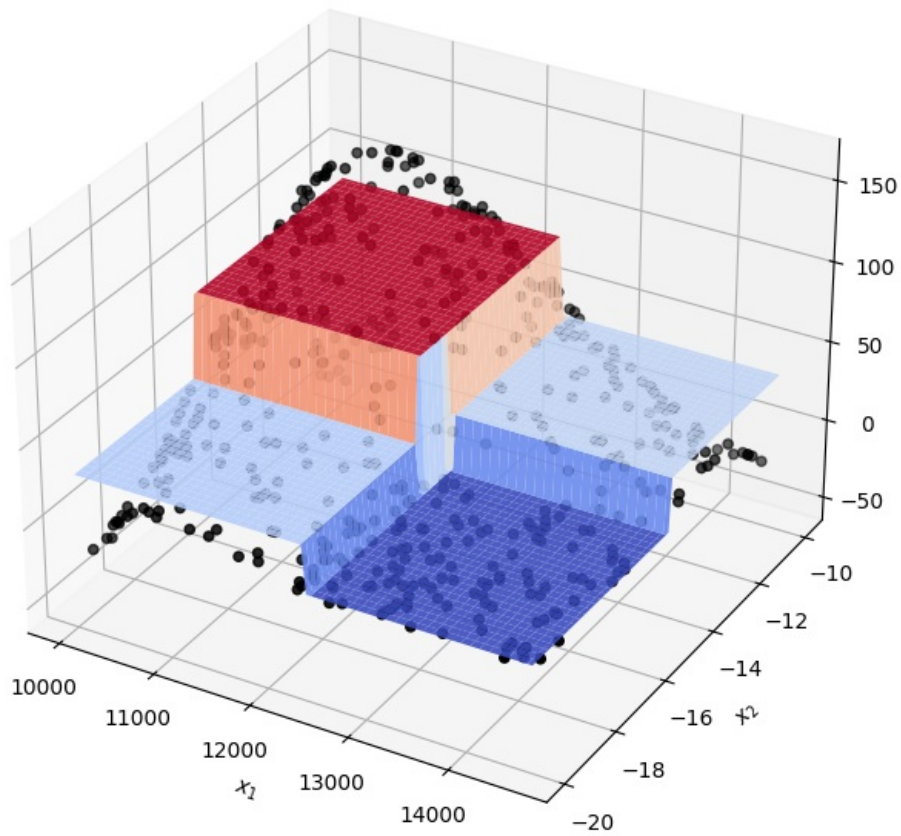
RMS error for Decision Tree Regressor on training data: 35.47184989095342

RMS error for Decision Tree Regressor on test data: 37.54886839401237

training data



test data

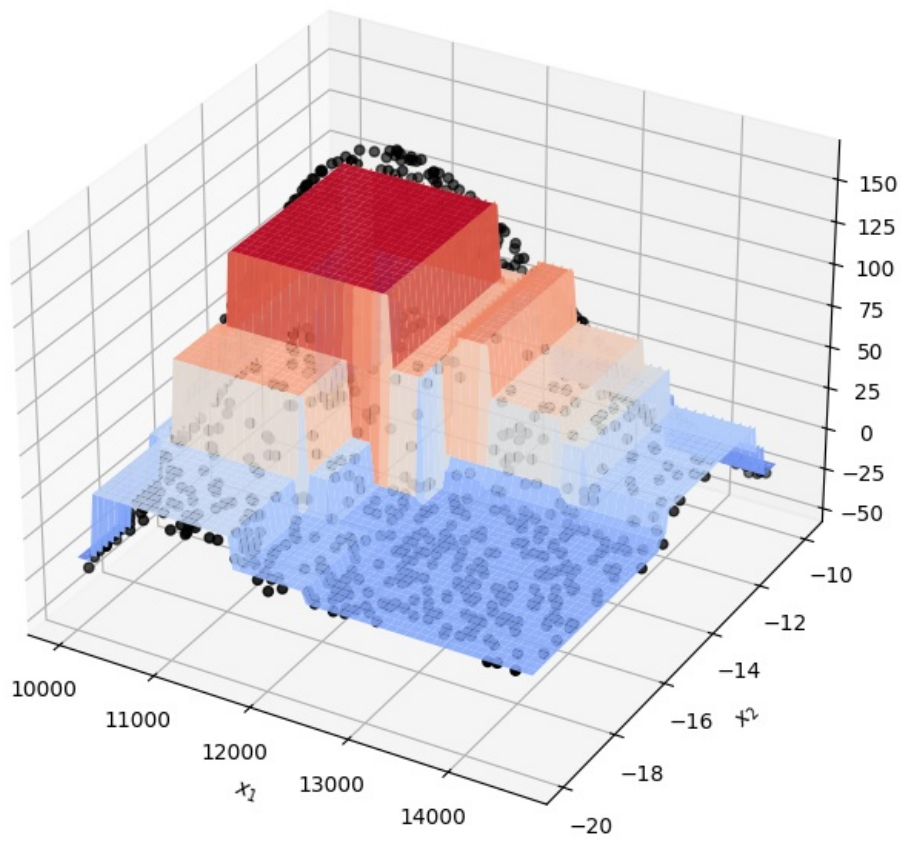


Max Depth: 5

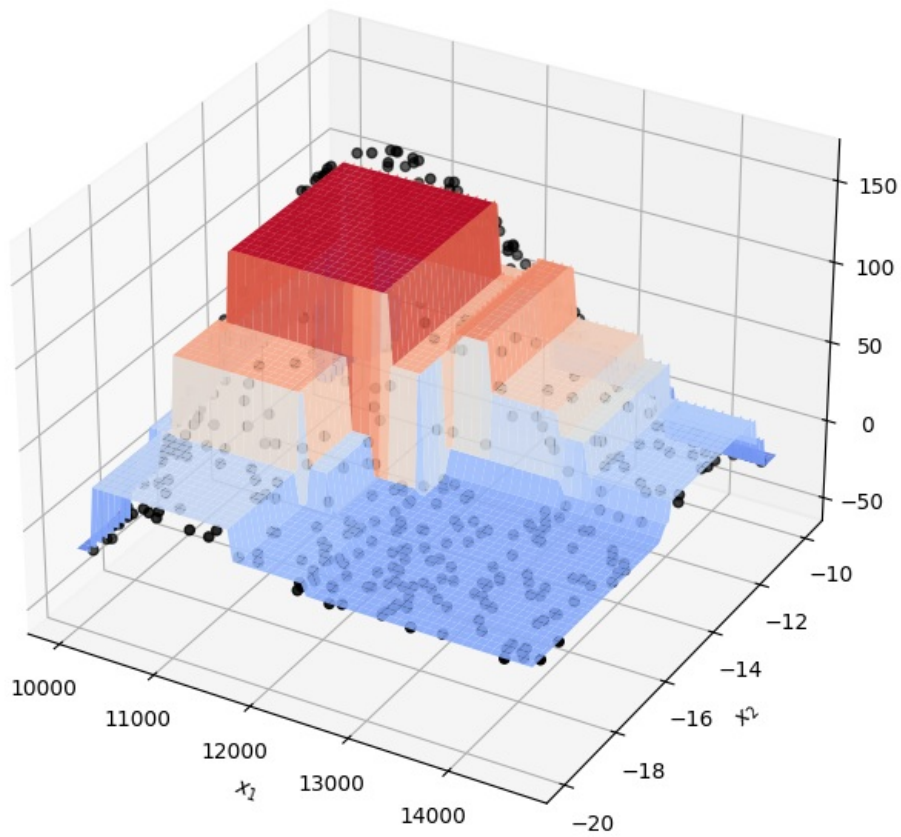
RMS error for Decision Tree Regressor on training data: 17.932673237502154

RMS error for Decision Tree Regressor on test data: 19.02935744931633

training data



test data

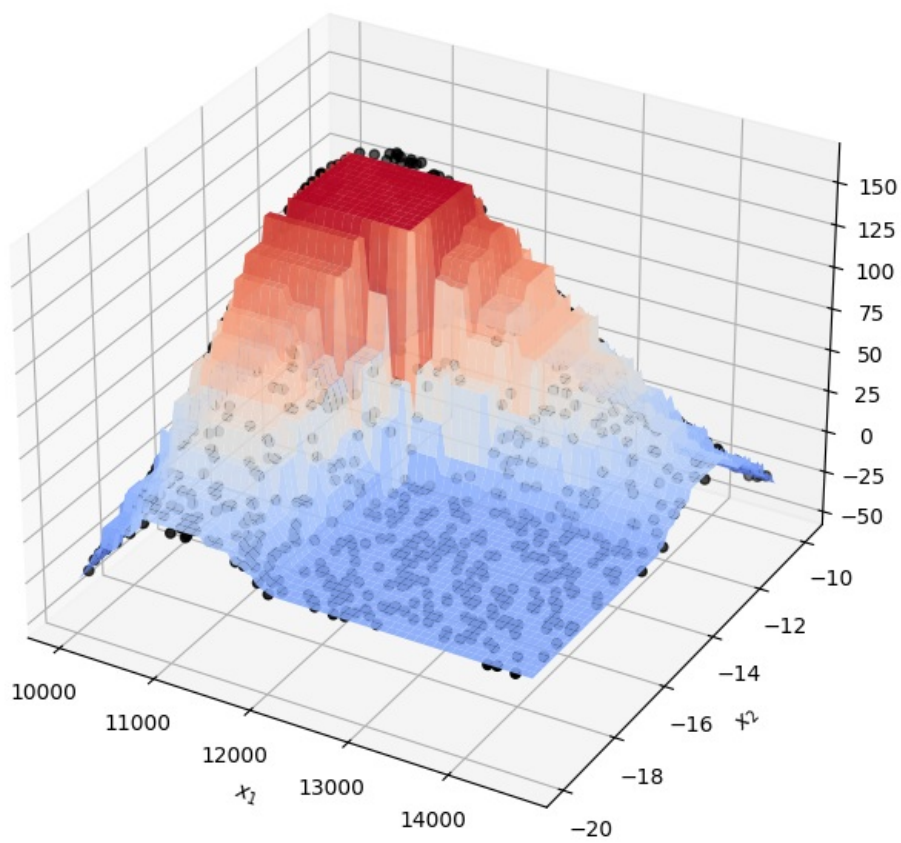


Max Depth: 10

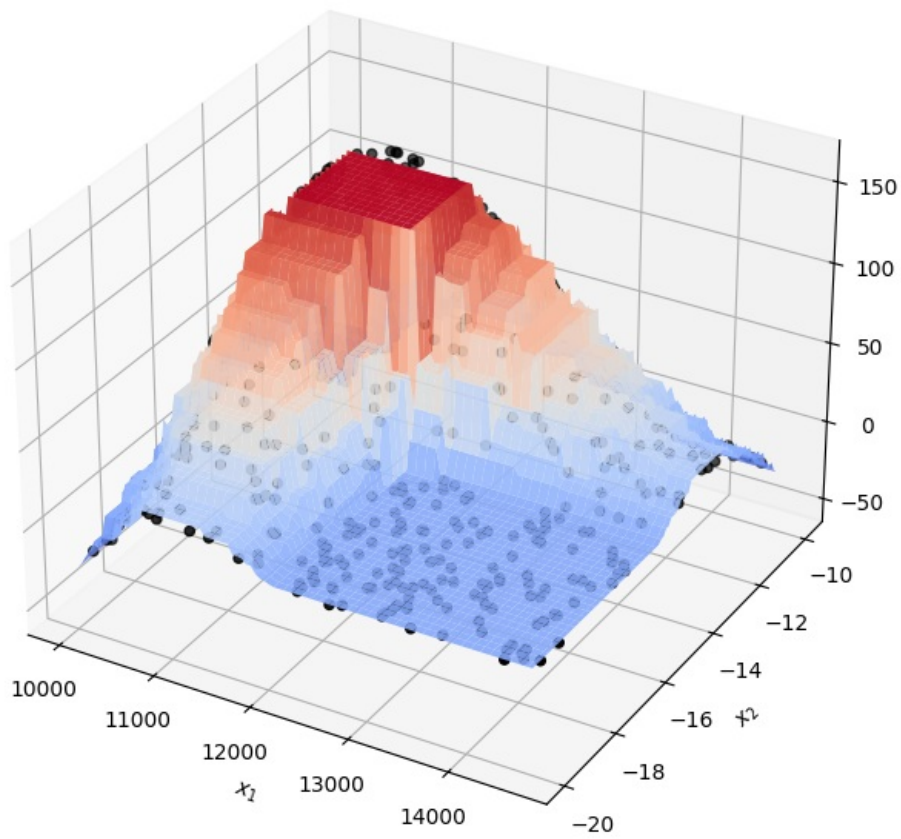
RMS error for Decision Tree Regressor on training data: 4.417134916147934

RMS error for Decision Tree Regressor on test data: 7.811633178930412

training data



test data

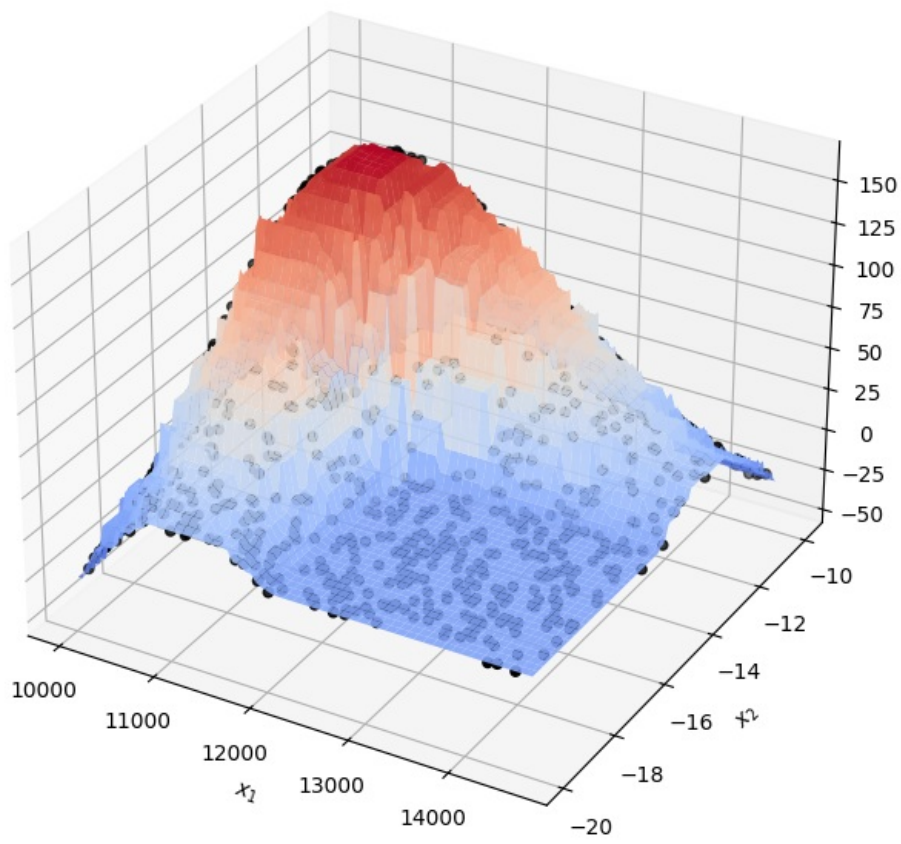


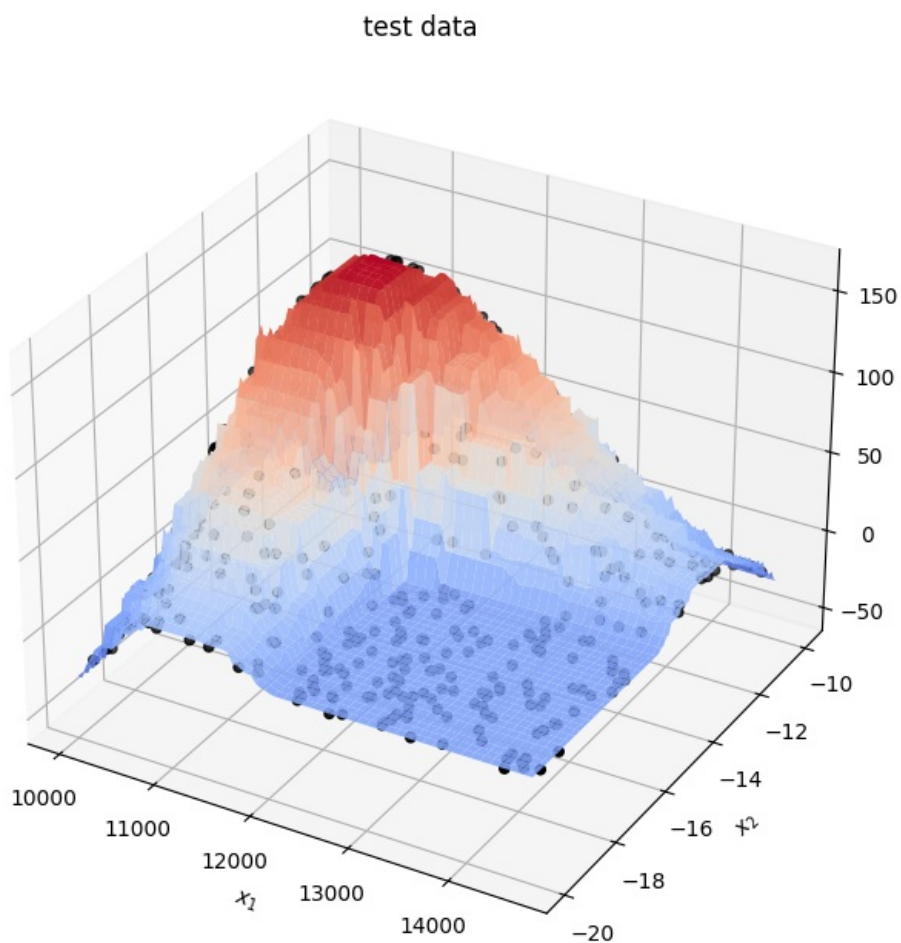
Max Depth: 25

RMS error for Decision Tree Regressor on training data: 0.0

RMS error for Decision Tree Regressor on test data: 5.597397150274175

training data





Question

- Which of your regression trees performed the best on testing data?
The model with the maximum depth of 25 performed best on the testing data.

Regression trees

Train 4 random forests in sklearn. For all of them, use the max depth values from your best-performing regression tree. The number of estimators should vary, with values [5, 10, 25, 100].

Plot the predictions as a surface plot along with test points. Once again, for each model, compute the train and test RMSE by calling your RMSE function. Print these results.

```
In [14]: # YOUR CODE GOES HERE
max_depth = [5, 10, 25, 100]

for i in max_depth:
```



```

dt = RandomForestRegressor(max_depth = i)
dt.fit(X_train, y_train)

pred_train = dt.predict(X_train)
pred_test = dt.predict(X_test)
rms_train = RMSE(y_train, pred_train)
rms_test = RMSE(y_test, pred_test)

print("\nMax Depth: ", i)
print("RMS error for Random Forest Regressor on training data: ", rms_train)
print("RMS error for Random Forest Regressor on test data: ", rms_test)
make_plot(X_train, y_train, dt, title = "training data")
make_plot(X_test, y_test, dt, title = "test data")

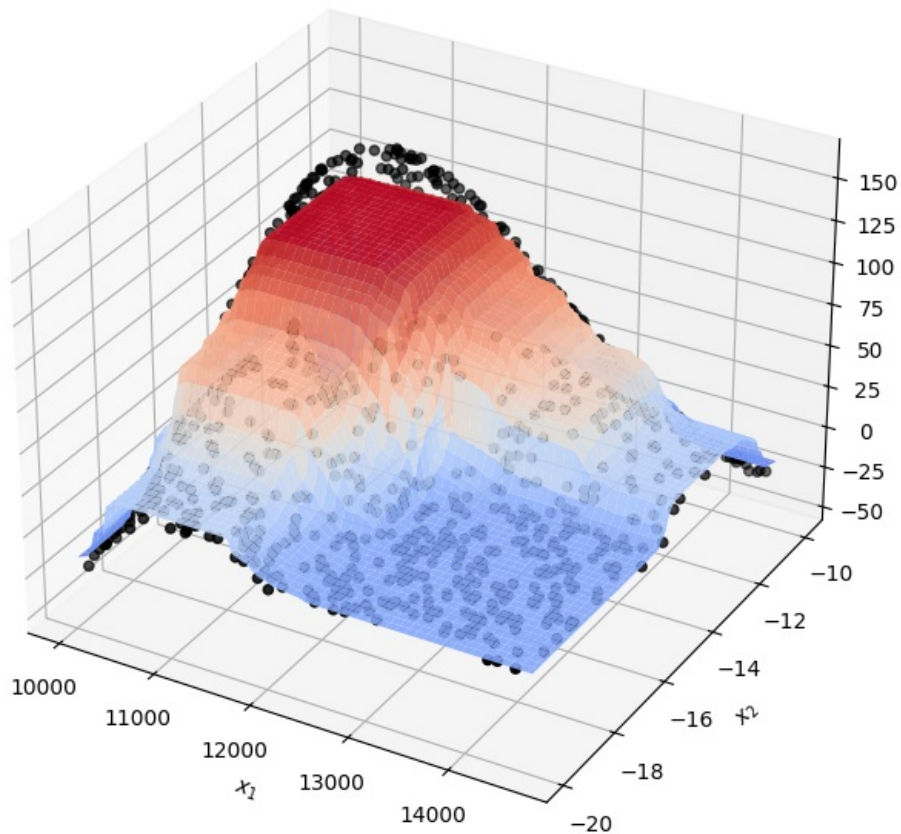
```

Max Depth: 5

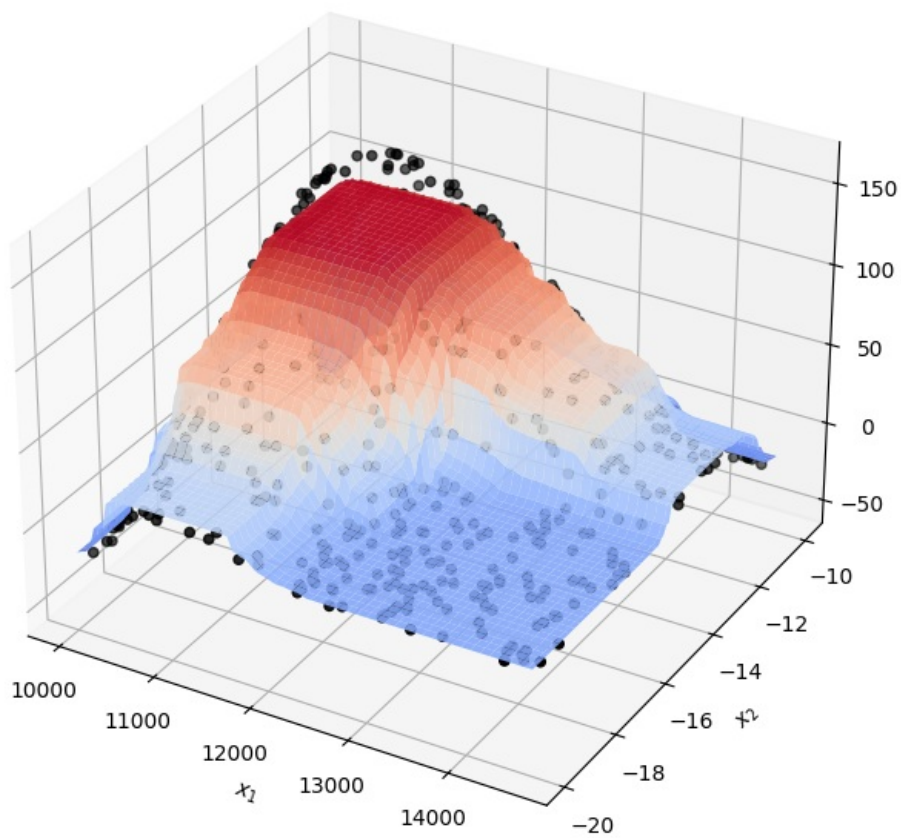
RMS error for Random Forest Regressor on training data: 14.870059581561991

RMS error for Random Forest Regressor on test data: 14.751953582837048

training data



test data

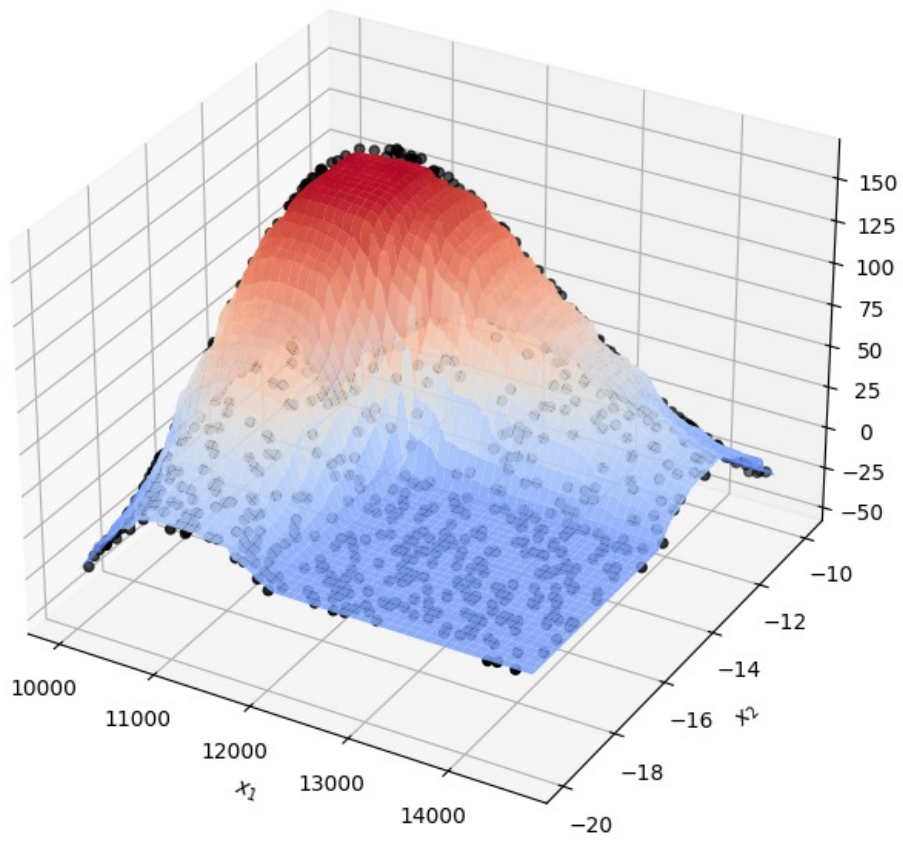


Max Depth: 10

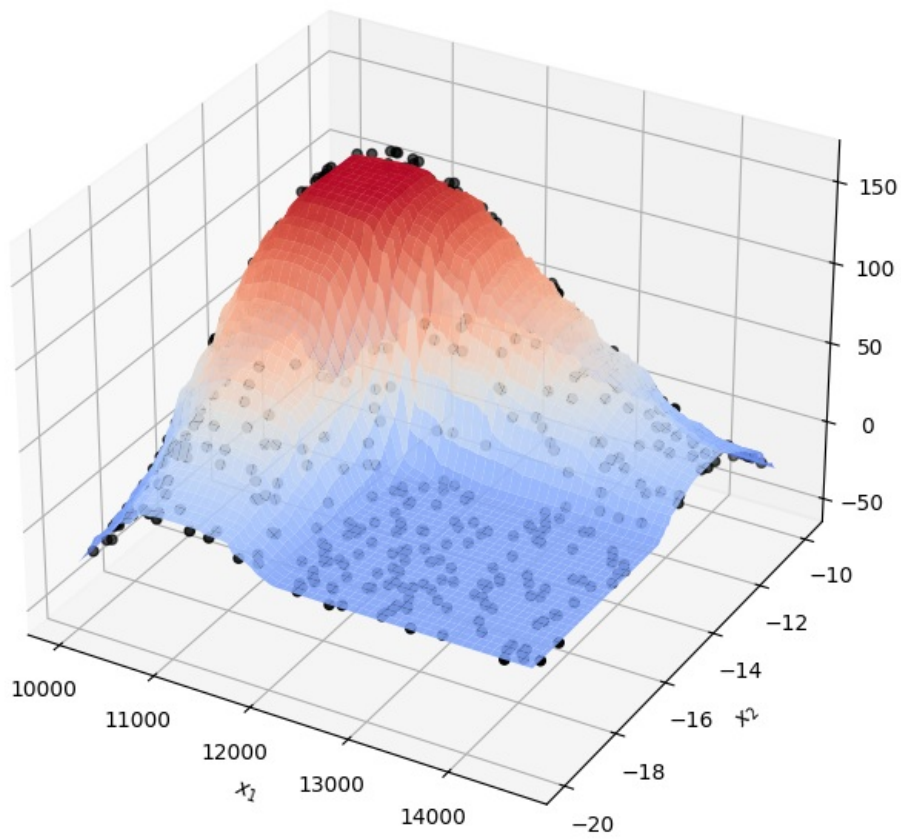
RMS error for Random Forest Regressor on training data: 2.946007563532823

RMS error for Random Forest Regressor on test data: 3.980096503464295

training data



test data

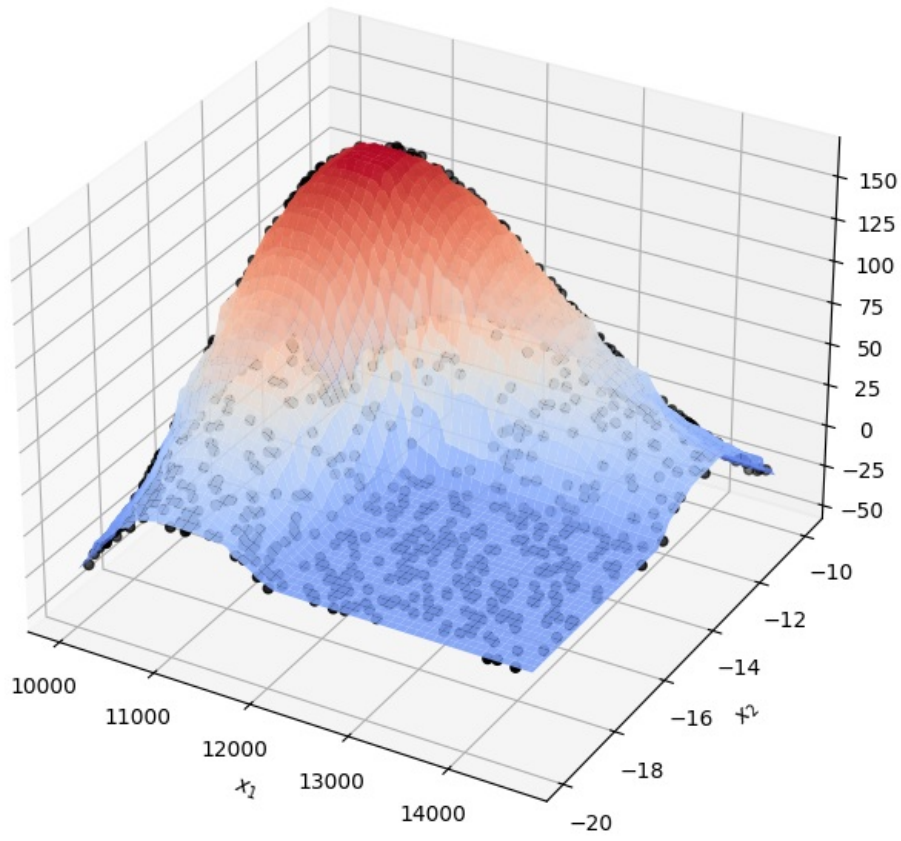


Max Depth: 25

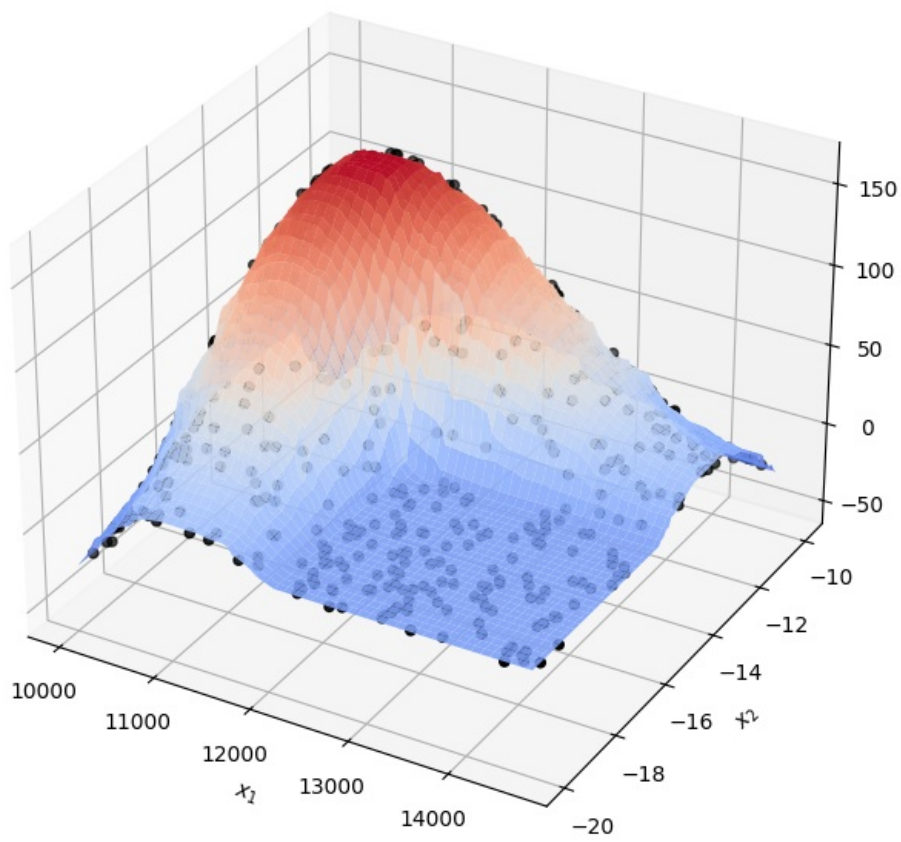
RMS error for Random Forest Regressor on training data: 1.3509754829019005

RMS error for Random Forest Regressor on test data: 3.0555287923736825

training data



test data

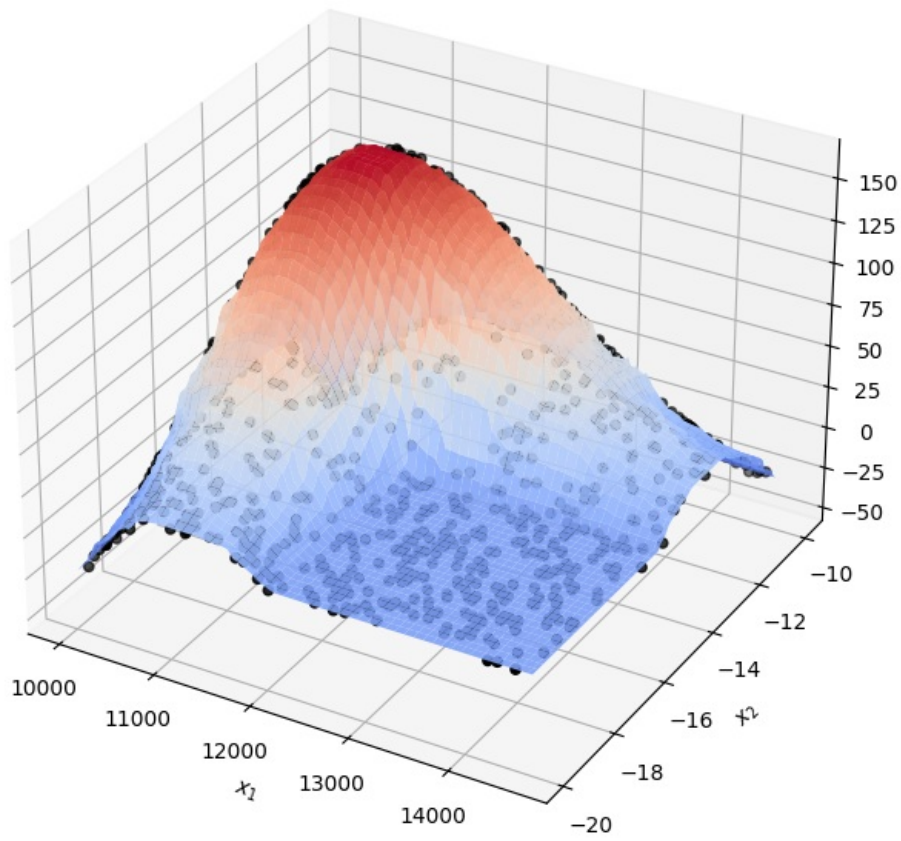


Max Depth: 100

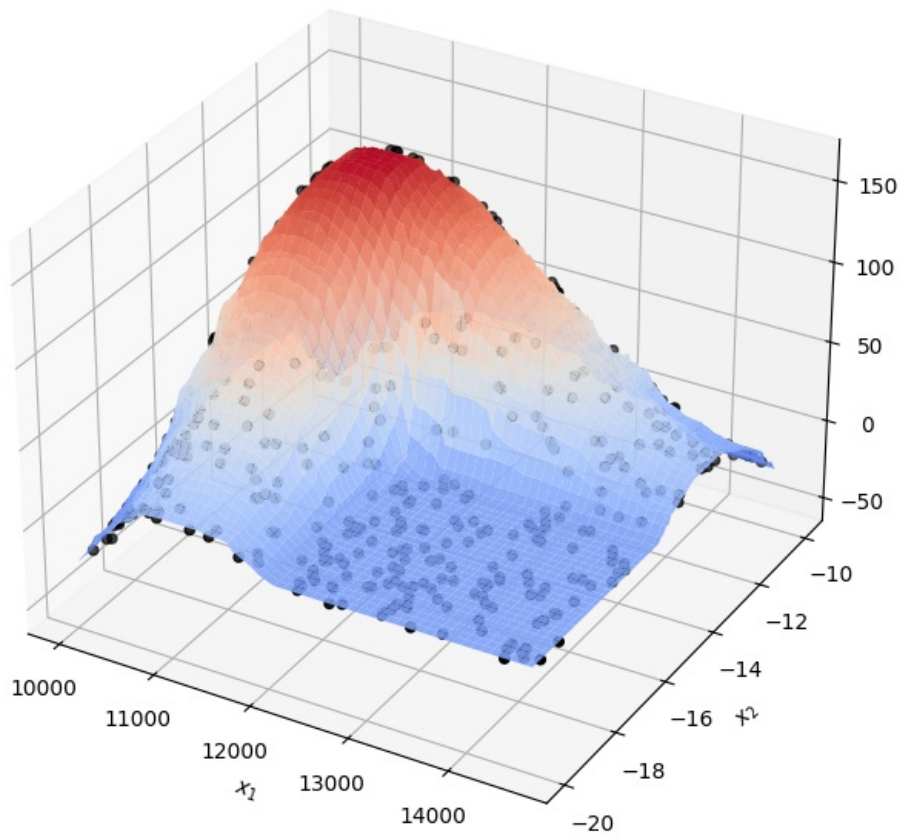
RMS error for Random Forest Regressor on training data: 1.3685013000875361

RMS error for Random Forest Regressor on test data: 2.976915452615326

training data



test data



Questions

- Which of your random forests performed the best on testing data?

Again, the model with the highest maximum depth performed the best on the testing data. This was slightly surprising, since I expected that the model with the highest maximum depth might would be overfit to the training data, resulting in sub-optimal performance in the test data.

- How does the random forest prediction surface differ qualitatively from that of the decision tree? The random forest prediction surface is much smoother and without irregularities when compared to that of the decision tree, which is much more discontinuous and jagged. Thanks to the averaging and randomization in the random forest model, leading to it being less sensitive to noise and outliers.

Loading [MathJax]/extensions/Safe.js


```
In [51]: from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.inspection import permutation_importance
import numpy as np
from sklearn.ensemble import RandomForestRegressor
import matplotlib.pyplot as plt

data = load_iris()
X = data.data
y = data.target
X_train, X_test, y_train, y_test = train_test_split(X, y, stratify = y, random_state = 10, test_size = 0.2)
```

```
In [52]: model = RandomForestRegressor(n_estimators = 5)
model.fit(X_train, y_train)
```

```
Out[52]: ▼      RandomForestRegressor
RandomForestRegressor(n_estimators=5)
```

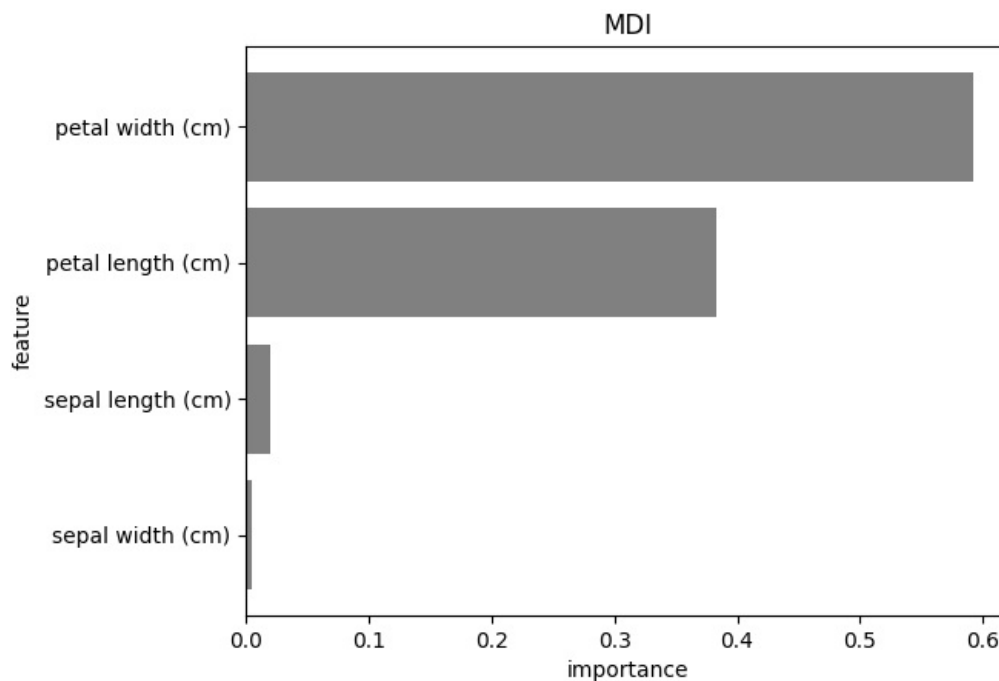
```
In [53]: mdi = model.feature_importances_
indices_mdi = np.argsort(mdi)

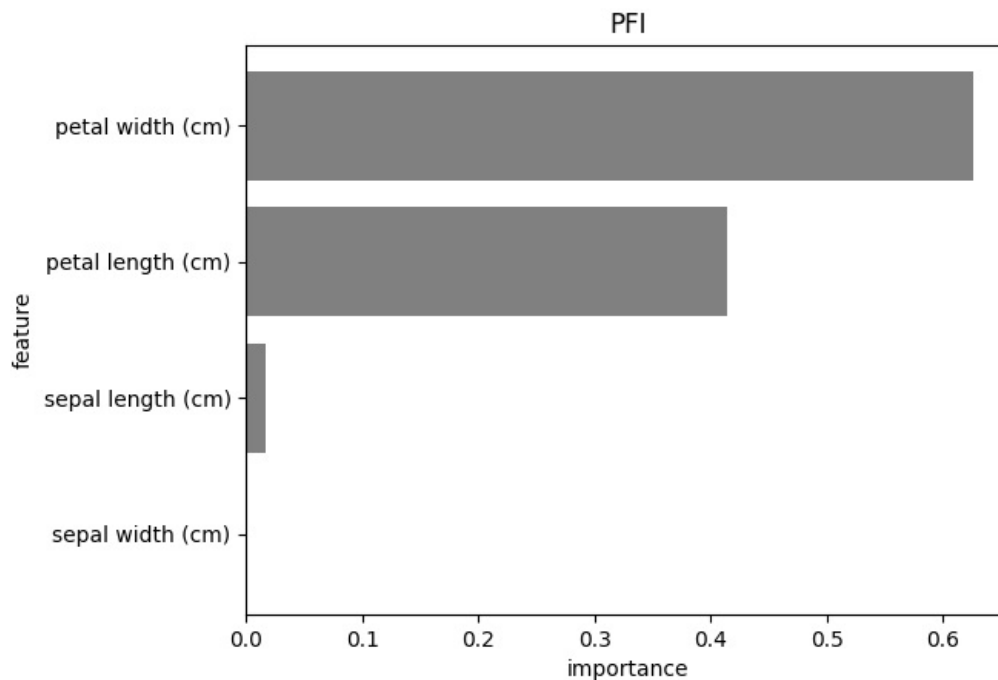
pfi_result = permutation_importance(model, X_test, y_test, n_repeats = 10, random_state = 50)
pfi = pfi_result.importances_mean
indices_pfi = np.argsort(pfi)
```

```
In [54]: labels = data.target_names
features = data.feature_names

plt.barh(np.array(features)[indices_mdi], np.array(mdi)[indices_mdi], color='gray')
plt.title("MDI")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()

plt.barh(np.array(features)[indices_pfi], np.array(pfi)[indices_pfi], color='gray')
plt.title("PFI")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()
```





```
In [55]: rand_feature = np.random.normal(loc = 2, scale = 1, size = (X.shape[0], 1))
X_ex = np.hstack([X, rand_feature])

X_train_ex, X_test_ex, y_train_ex, y_test_ex = train_test_split(X_ex, y, stratify = y, random_state = 10, test_size = 0.2)

model_ = RandomForestRegressor(n_estimators = 5)
model_.fit(X_train_ex, y_train_ex)
```

```
Out[55]: Random Forest Regressor
RandomForestRegressor(n_estimators=5)
```

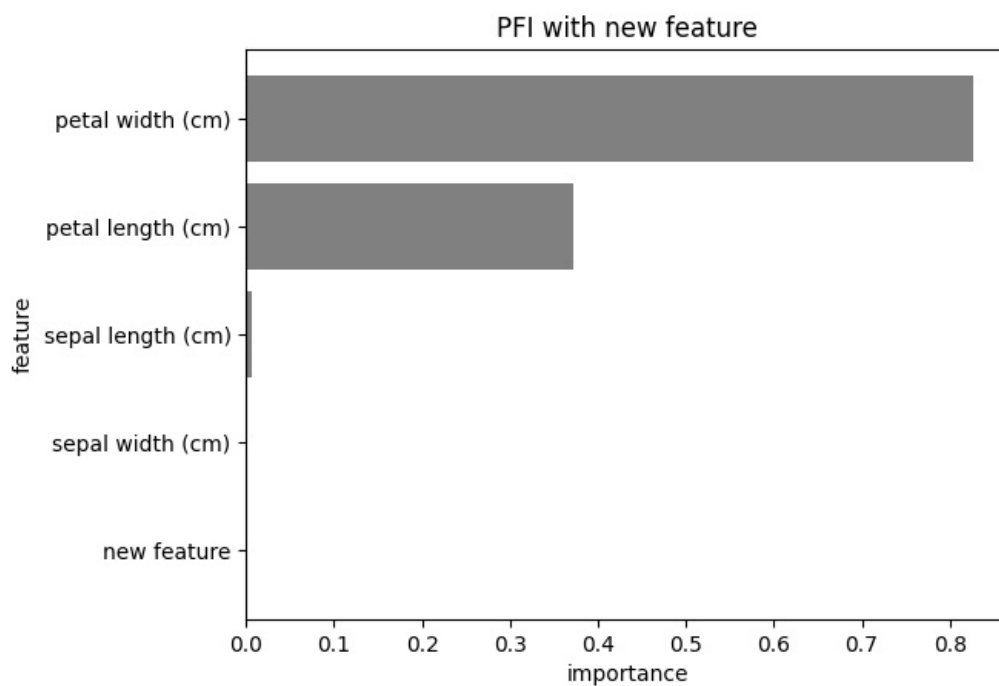
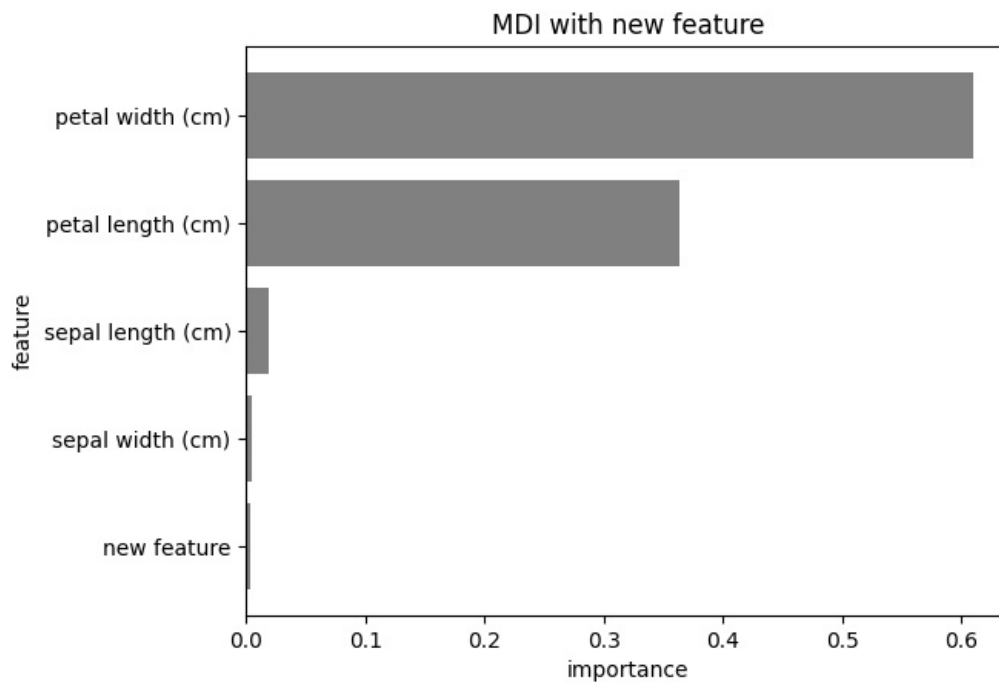
```
In [56]: mdi_ = model_.feature_importances_
indices_mdi_ = np.argsort(mdi_)

pfi_result_ = permutation_importance(model_, X_test_ex, y_test_ex, n_repeats = 10, random_state = 50)
pfi_ = pfi_result_.importances_mean
indices_pfi_ = np.argsort(pfi_)
```

```
In [57]: features_ = features
features_.append("new feature")

plt.barh(np.array(features_)[indices_mdi_], np.array(mdi_)[indices_mdi_], color='gray')
plt.title("MDI with new feature")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()

plt.barh(np.array(features_)[indices_pfi_], np.array(pfi_)[indices_pfi_], color='gray')
plt.title("PFI with new feature")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()
```



```
In [58]: X_exx = np.hstack([X_ex, np.ones([X.shape[0], 1])])
X_train_exx, X_test_exx, y_train_exx, y_test_exx = train_test_split(X_exx, y, stratify = y, random_state = 10,
model_1 = RandomForestRegressor(n_estimators = 5)
model_1.fit(X_train_exx, y_train_exx)
```

```
Out[58]: Random Forest Regressor
RandomForestRegressor(n_estimators=5)
```

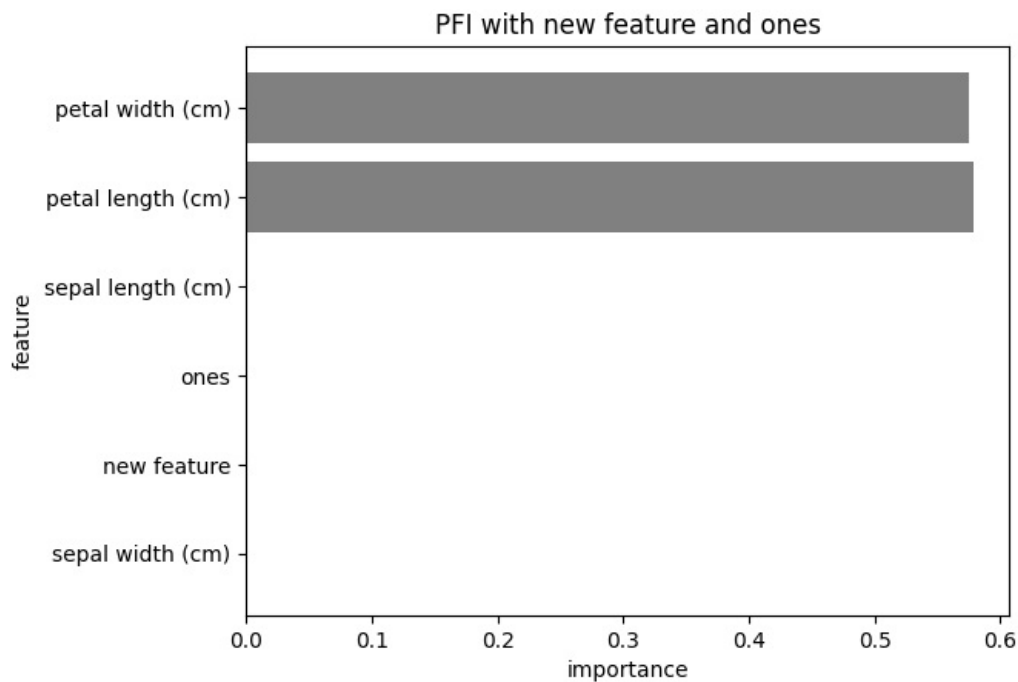
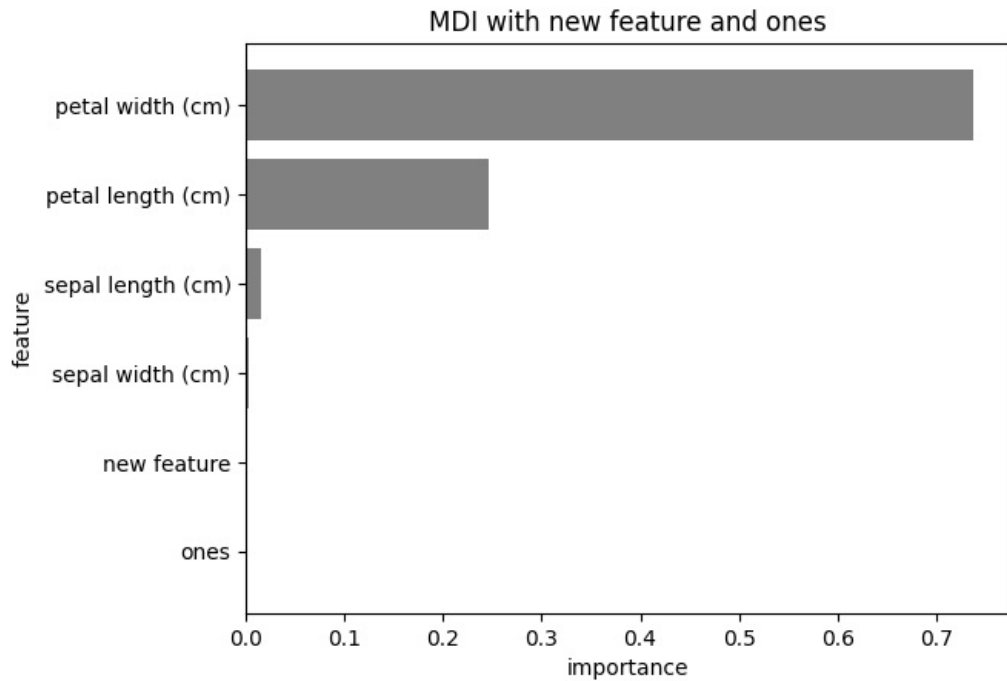
```
In [59]: mdi_1 = model_1.feature_importances_
indices_mdi_1 = np.argsort(mdi_1)

pfi_result_1 = permutation_importance(model_1, X_test_exx, y_test_exx, n_repeats = 10, random_state = 50)
pfi_1 = pfi_result_1.importances_mean
indices_pfi_1 = np.argsort(pfi_1)
```

```
In [60]: features_1 = features_
features_1.append("ones")

plt.barh(np.array(features_1)[indices_mdi_1], np.array(mdi_1)[indices_mdi_1], color = 'gray')
plt.title("MDI with new feature and ones")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()
```

```
plt.barh(np.array(features_1)[indices_pfi_1], np.array(pfi_1)[indices_pfi_1], color = 'gray')
plt.title("PFI with new feature and ones")
plt.ylabel("feature")
plt.xlabel("importance")
plt.show()
```



Conclusions

1. MDI might rank random or irrelevant features as more important than they actually are because it depends on how frequently a feature is used to split the data, not whether it improves the model's predictive performance.
2. PFI is a more robust measure for identifying features that really contribute to the model's performance, as it evaluates how sensitive the model's accuracy is to shuffling the feature values.
3. The constant feature is correctly identified by both MDI and PFI as unimportant, but MDI can be tricked into ranking random features higher than they deserve.