

Lab 6: Decision Trees and Random Forests

In this lab you'll train a decision tree classifier and a random forest classifier. You'll do so on both synthetic and real data.

Run the code cell below to import the required packages.

In [1]:

```
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
import sklearn
import sklearn.tree # For DecisionTreeClassifier class
import sklearn.ensemble # For RandomForestClassifier class
import sklearn.datasets # For make_circles
import sklearn.metrics # For accuracy_score
```

1. Fitting a Random Forest classifier to synthetic data

Exercises 1.1-1.4 ask you to apply scikit-learn's decision tree classifier ([sklearn.tree.DecisionTreeClassifier](#)) and random forest classifier ([sklearn.ensemble.RandomForestClassifier](#)) to synthetic data.

Run the code cell below to define some useful functions for plotting data and predictions.

In [2]:

```
def plot_data(X, y):
    """Plots a toy 2D data set. Assumes values in range [-3,3] and at most 3 classes."""
    plt.plot(X[y==0,0], X[y==0,1], 'ro', markersize=6)
    plt.plot(X[y==1,0], X[y==1,1], 'bo', markersize=6)
    X = np.column_stack([x1.ravel(), x2.ravel(), 0])
    y = model.predict(X).reshape(x1.shape)
    plt.xlim(-3, 3)
    plt.ylim(-3, 3)
    plt.xlabel('x1')
    plt.ylabel('x2')
    plt.gca().set_aspect('equal')

def plot_predict(model):
    """Plots the model's predictions over all points in range 2D [-3, 3].
    Assumes at most 3 classes.
    """
    extent = (-3, 3, -3, 3)
    xmin, xmax, x2min, x2max = extent
    x1, x2 = np.meshgrid(np.linspace(xmin, xmax, 100), np.linspace(x2min, x2max, 100))
    X = np.column_stack([x1.ravel(), x2.ravel(), 0])
    cmap = matplotlib.colors.ListedColormap(['r', 'b', 'g'])
    plt.imshow(model.predict, origin='lower', alpha=0.4, vmin=0, vmax=2, cmap=cmap, interpolation='nearest')
    plt.xlim(x2min, x2max)
    plt.ylim(x1min, x1max)
    plt.gca().set_aspect('equal')

def plot_class_probability(model, class_index):
    """Plots the model's class probability for the given class (0,1,2)
    over all points in range 2D [-3, 3]. Assumes at most 3 classes.
    """
    extent = (-3, 3, -3, 3)
    xmin, xmax, x2min, x2max = extent
    x1, x2 = np.meshgrid(np.linspace(xmin, xmax, 100), np.linspace(x2min, x2max, 100))
    X = np.column_stack([x1.ravel(), x2.ravel(), 0])
    p = model.predict_proba(X)[:, class_index].reshape(x1.shape)
    colors = [(1, 0, 0), (0, 1, 1), (0, 1, 1), (0, 1, 1), (0, 1, 1)]
    cmap = matplotlib.colors.ListedColormap(np.linspace(1, 1, 1), colors[class_index], 50)
    plt.imshow(p, extent=extent, origin='lower', alpha=0.4, vmin=0, vmax=1, cmap=cmap)
    plt.xlim(x2min, x2max)
    plt.ylim(x1min, x1max)
    plt.gca().set_aspect('equal')
```

Exercise 1.1 — Train and inspect a small decision tree (2 points, 2 classes)

Read the documentation for the [DecisionTreeClassifier](#)'s `fit` method. Notice that the `y` vector should contain integer class labels. You are asked to build the small 2D training set below:

$$X = \begin{bmatrix} -1 & 0 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Write a few lines of code to

1. Define the training set above in two variables `X` and `y`.
2. Train a decision tree classifier on `X` and `y`. Use argument `random_state=0`.
3. Plot the decision tree predictions and the data (use `plot_predict` and `plot_data` from preamble).
4. Plot the decision tree itself (use `sklearn.tree.plot_tree`); pass `feature_names=['x1', 'x2']` as an argument.

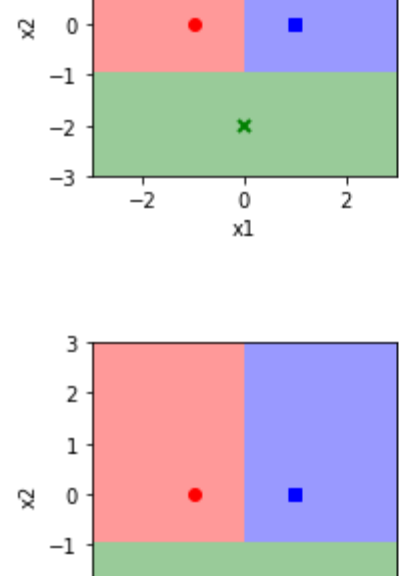
You should end up with a plot showing the data and the trained decision surface between classes 0 (red) and class 1 (blue). You should see a binary decision tree diagram depicting a tree of height 1 that splits the feature space using the first variable (`x1`) at threshold 0.

Tip 1: If you want a single code cell to generate multiple figures, call `plt.figure()` to tell Matplotlib that you want the subsequent plotting commands to generate a new, separate figure from any previous plotting commands.

Tip 2: If the last line of your code returns a value, it will be printed as the `Out:` of the cell before the plots are shown. Some times you don't care about this 'final' value, for example if it is a string or some object you don't need/need. If you want to suppress the cell's `Out` value, add a semicolon (`;`) to the end of the last line of code in the cell.

In [13]:

```
# Your code here. Aim for 7-9 lines.
X = np.array([[ -1, 0], [ 1, 0], [ -1/3, 0], [ 1/3, 0]])
y = np.array([ 1, 1, 1, 1, 0])
decisionTree = sklearn.tree.DecisionTreeClassifier(random_state=0).fit(X, y)
plot_data(X, y)
plot_predict(decisionTree)
plt.figure()
sklearn.tree.plot_tree(decisionTree, feature_names=['x1', 'x2'])
```



Once you have things working, **add two data points** to your training set:

- $x_3 = (-\frac{1}{3}, 0)$ with class label $y_3 = 1$ (blue), and
- $x_4 = (\frac{1}{3}, 0)$ with class label $y_4 = 0$ (red).

Re-run your code cell above and make sure you understand how the splits and thresholds you see in the tree correspond to the decision region shown.

Note: If a decision tree node is shown as having values `[1,2]`, it means that node's region (between splitting) contains exactly one training point from class 0 and two training points from class 1. The root node thus 'contains' the entire training set.

Exercise 1.2 — Train and inspect a small decision tree (3 points, 3 classes)

Repeat Exercise 1.1 but with the following changes:

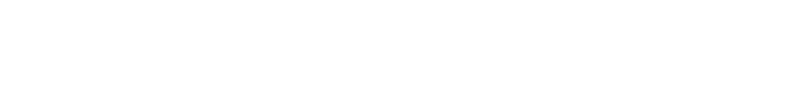
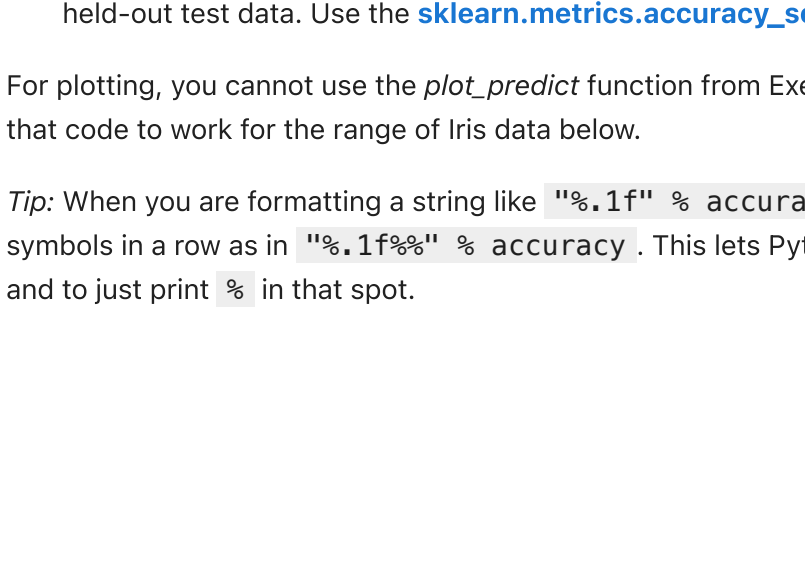
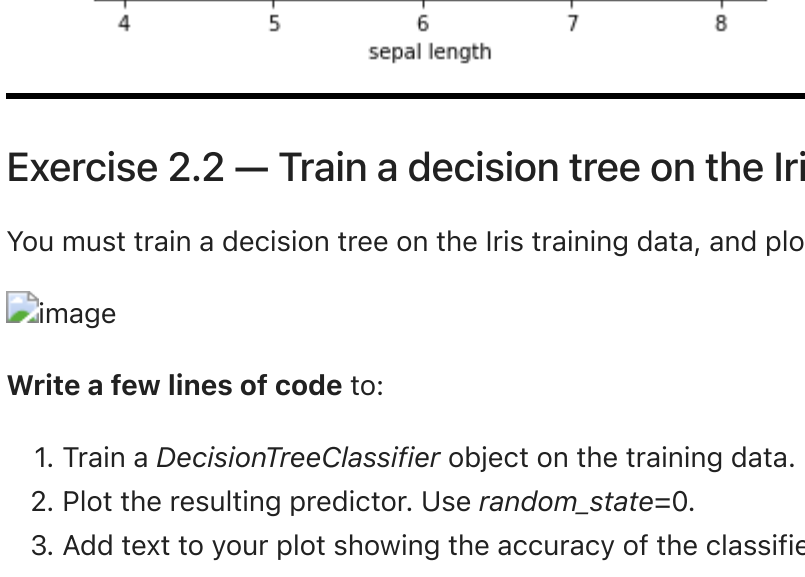
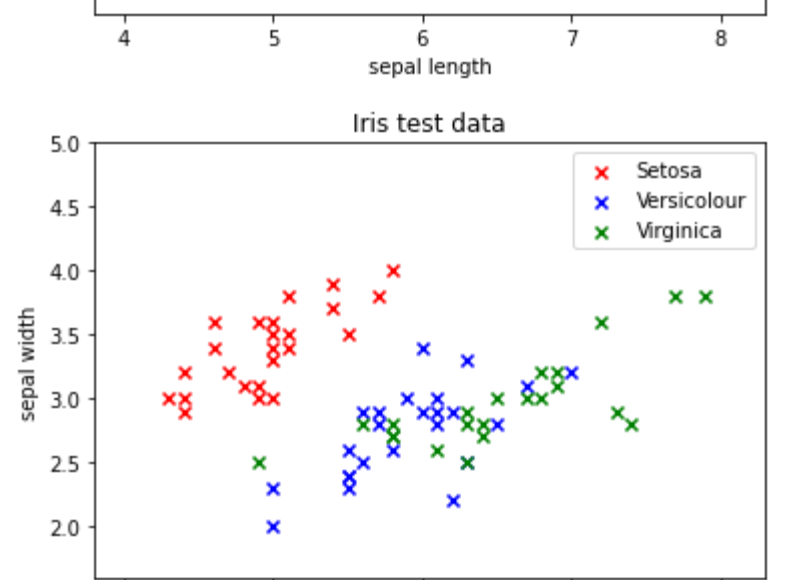
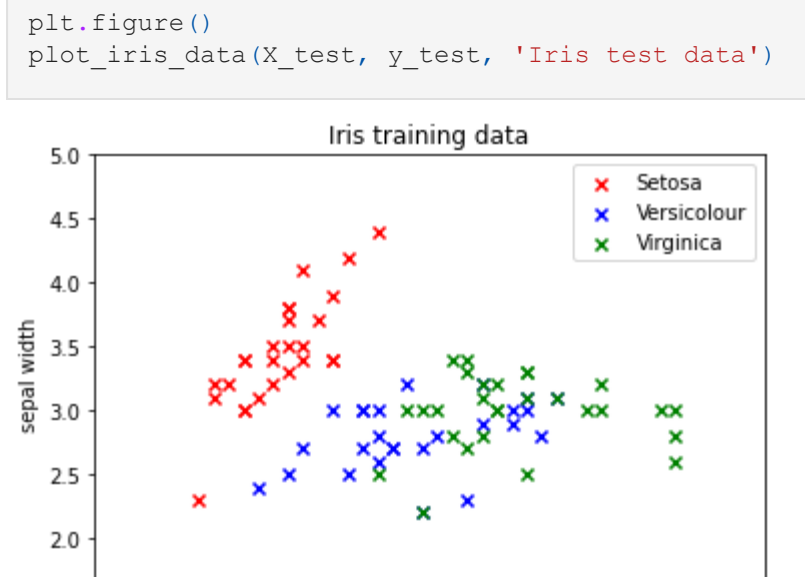
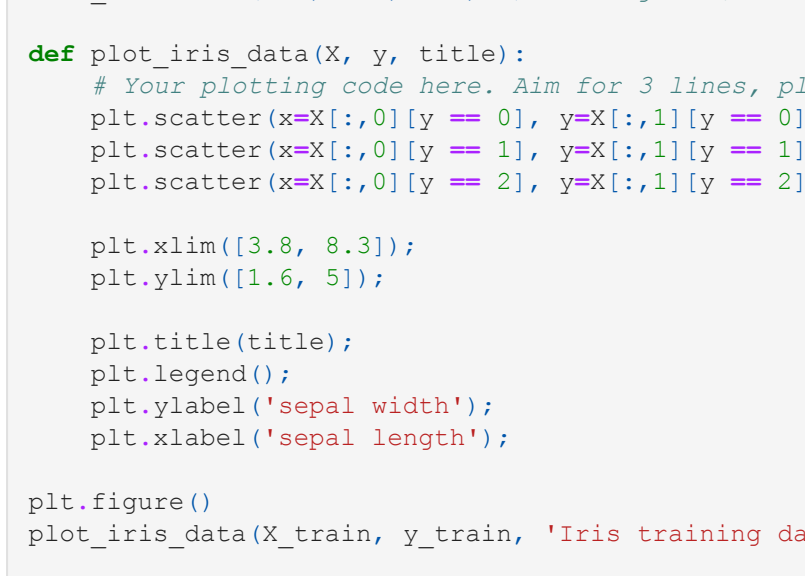
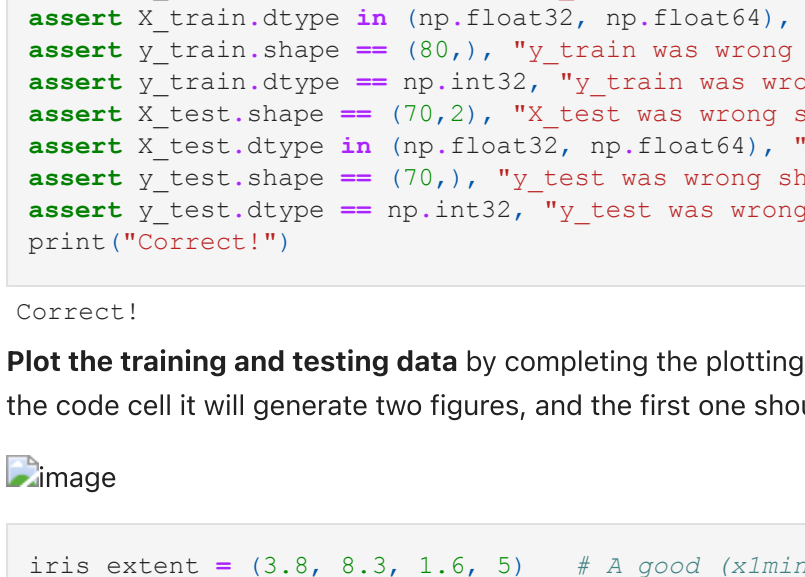
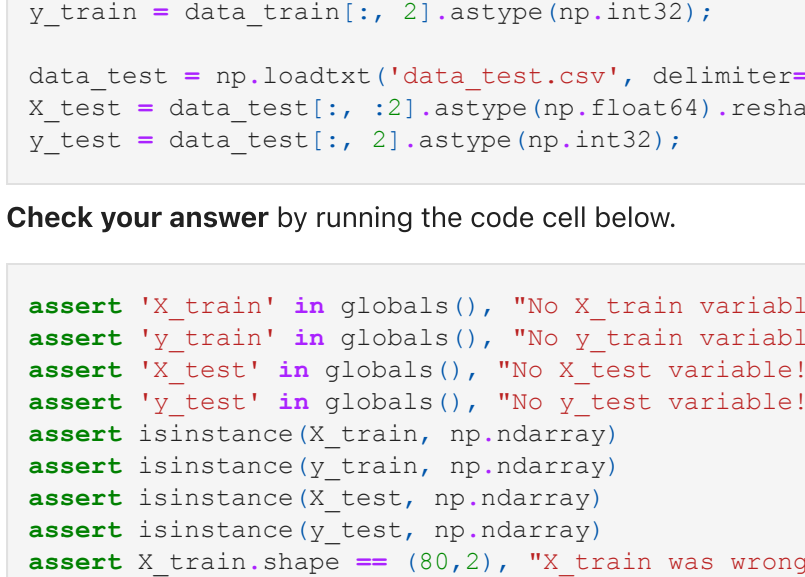
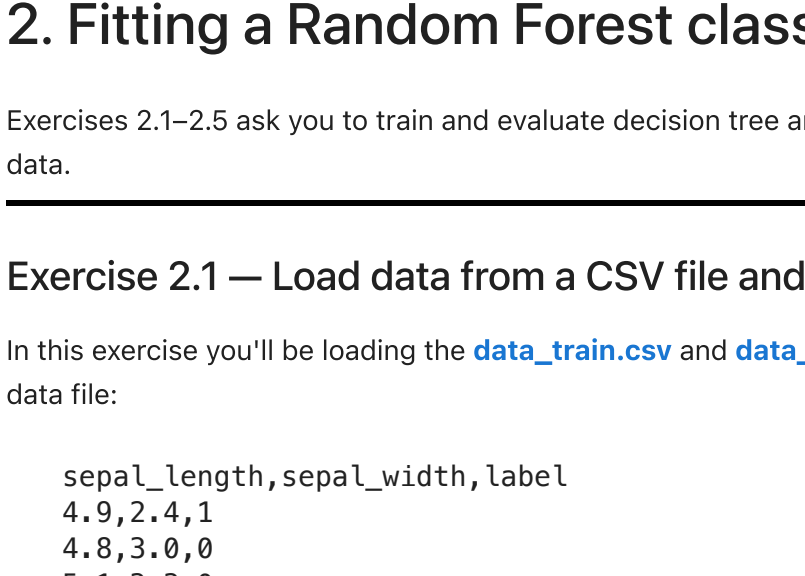
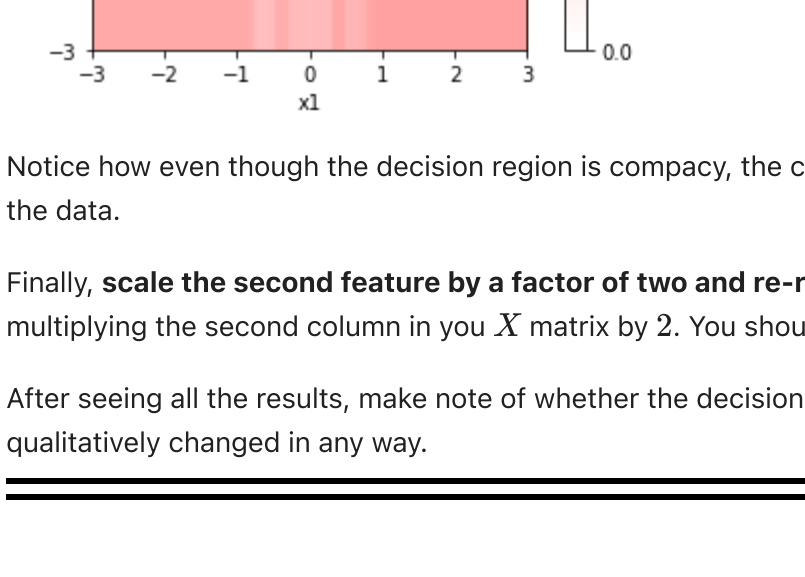
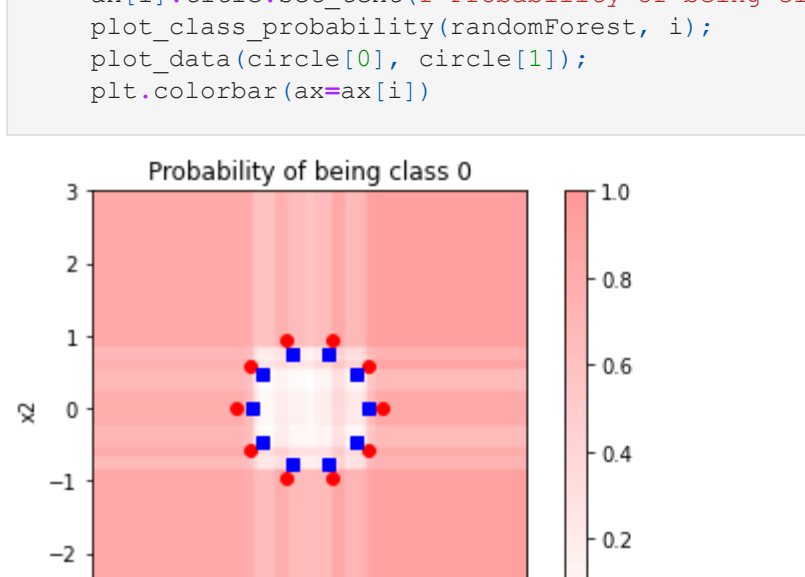
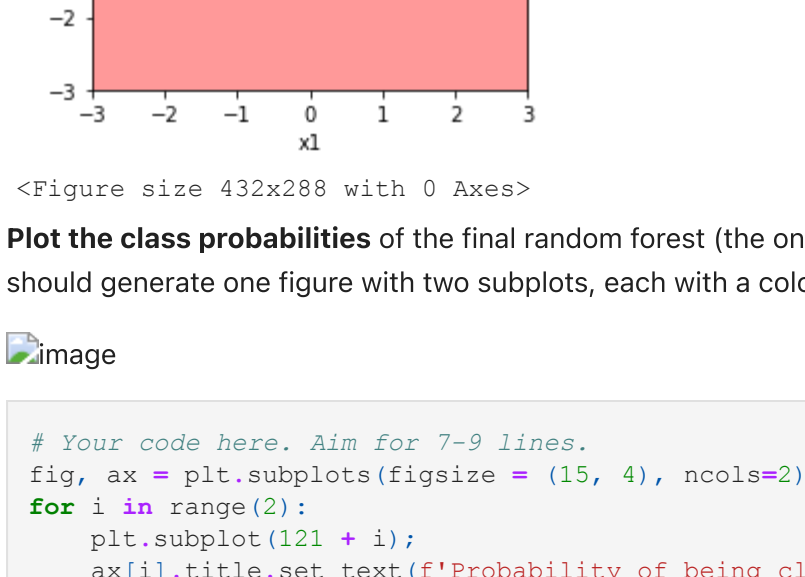
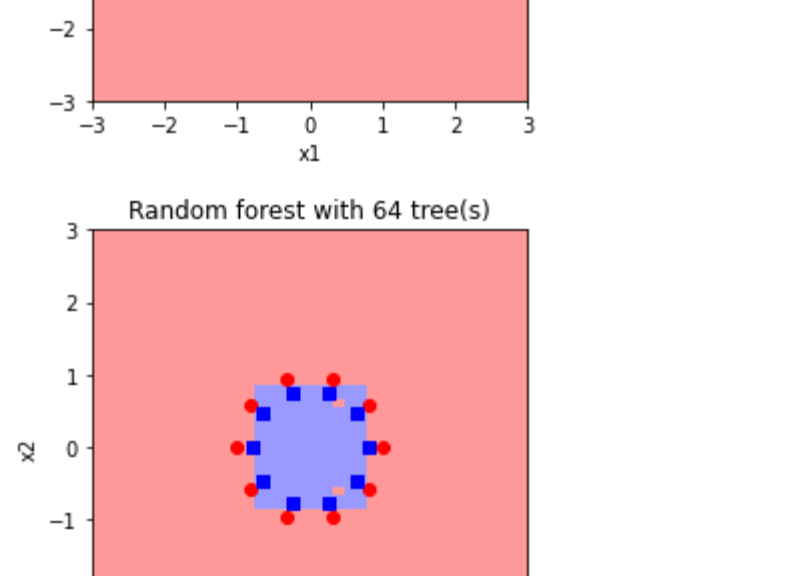
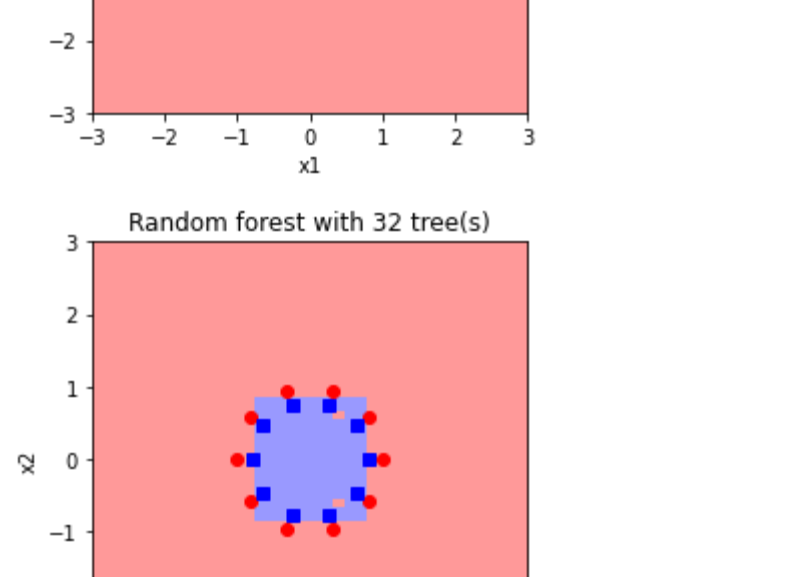
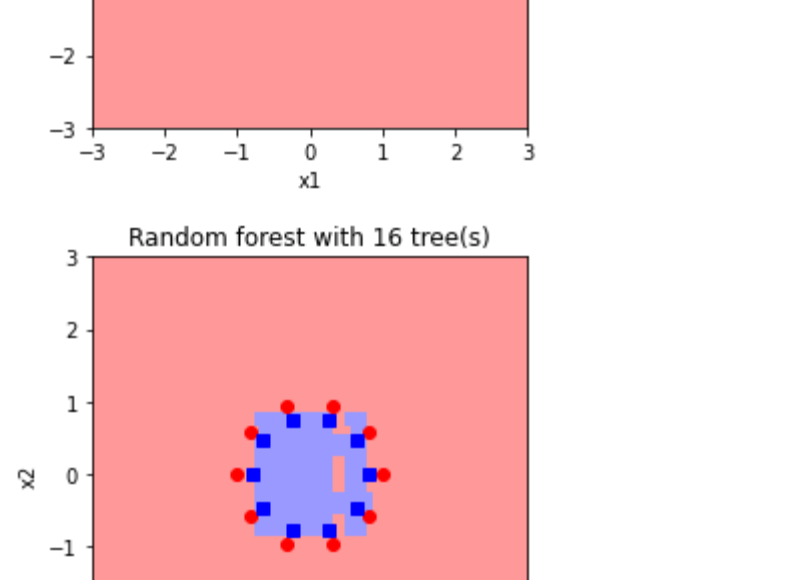
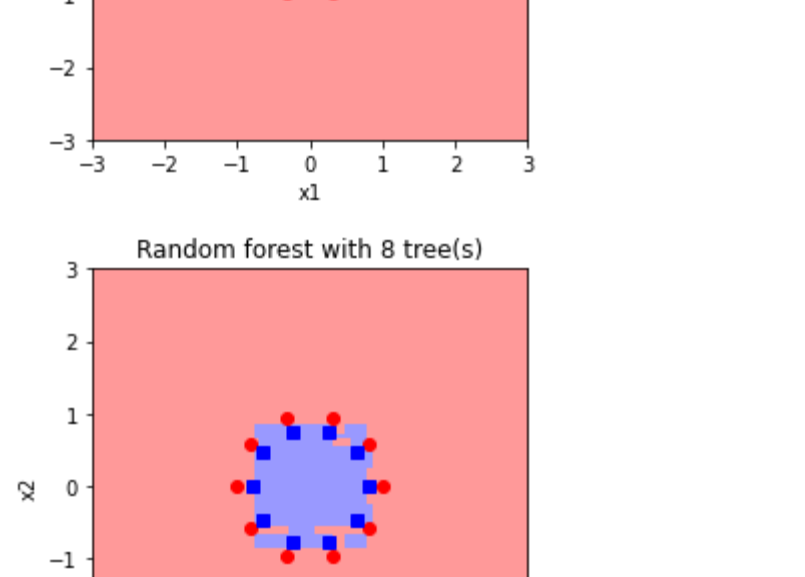
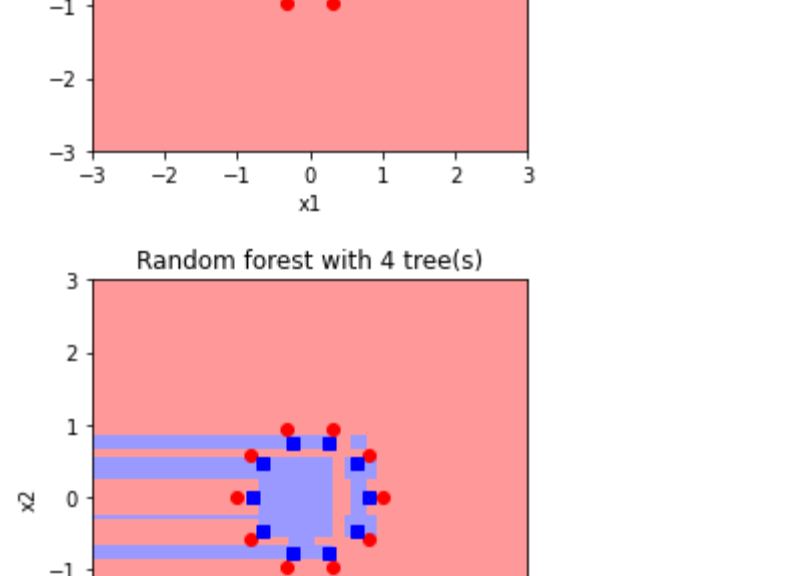
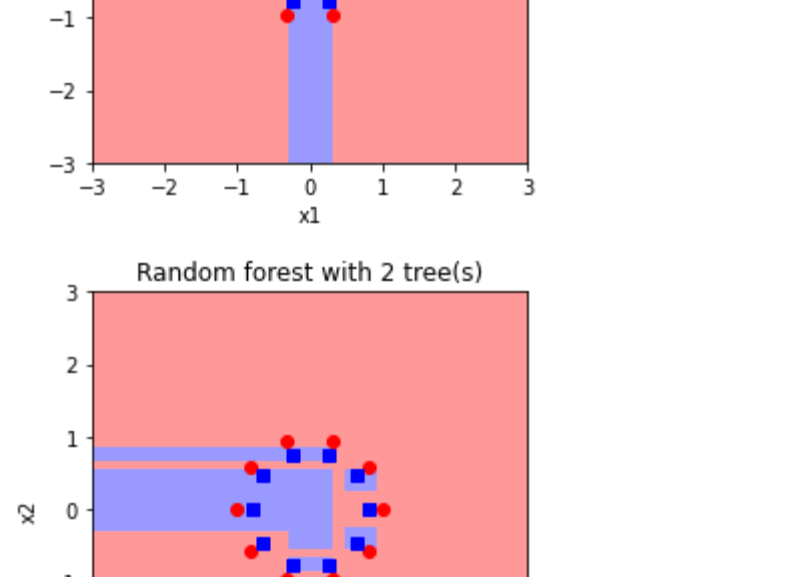
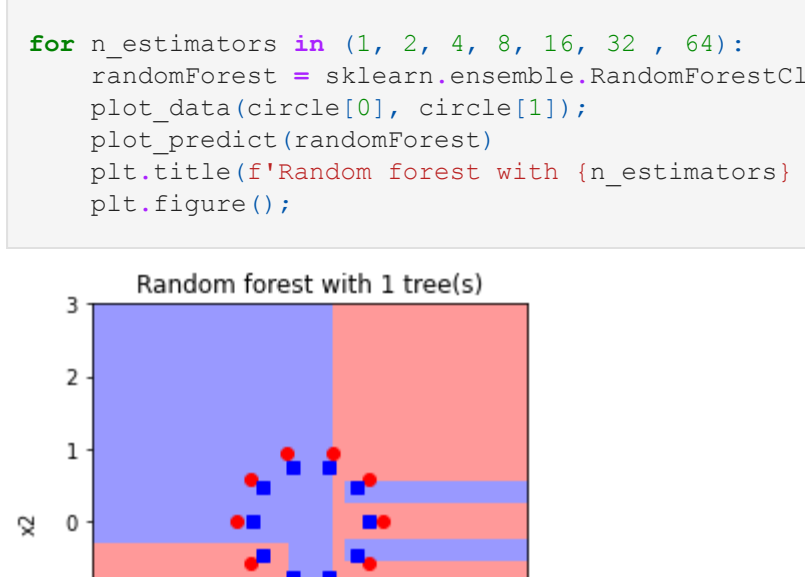
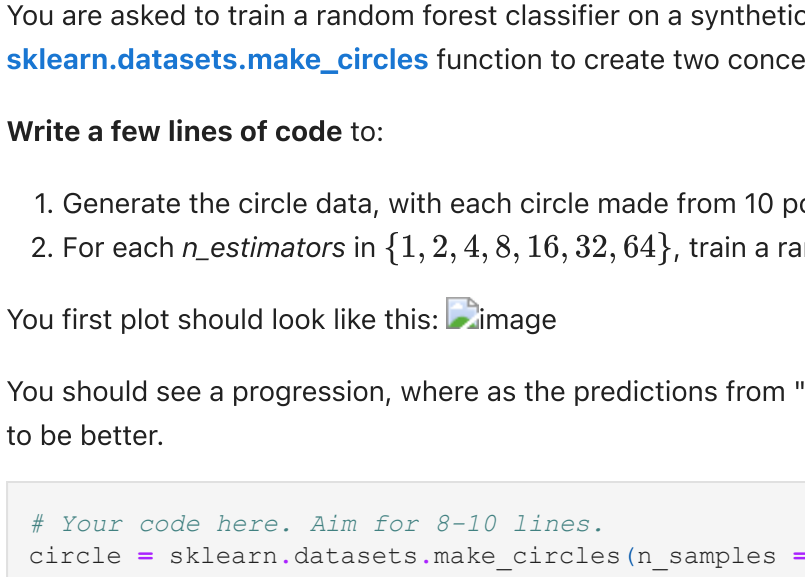
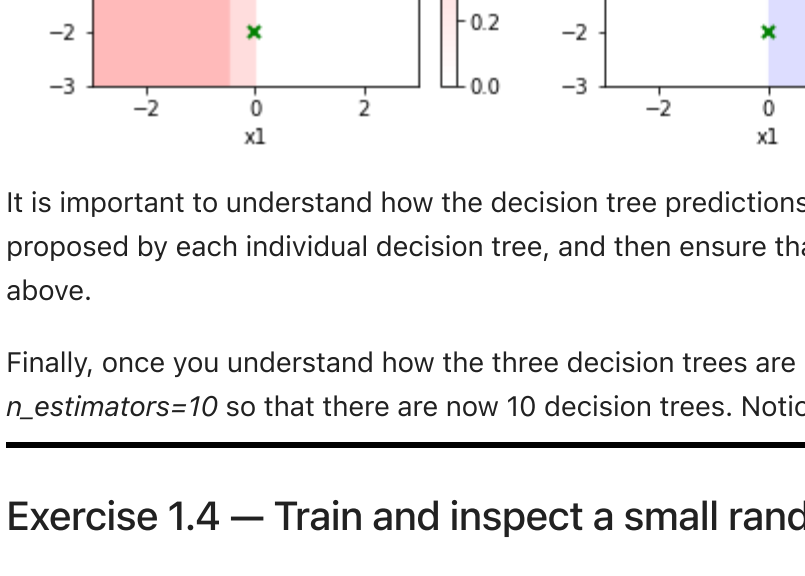
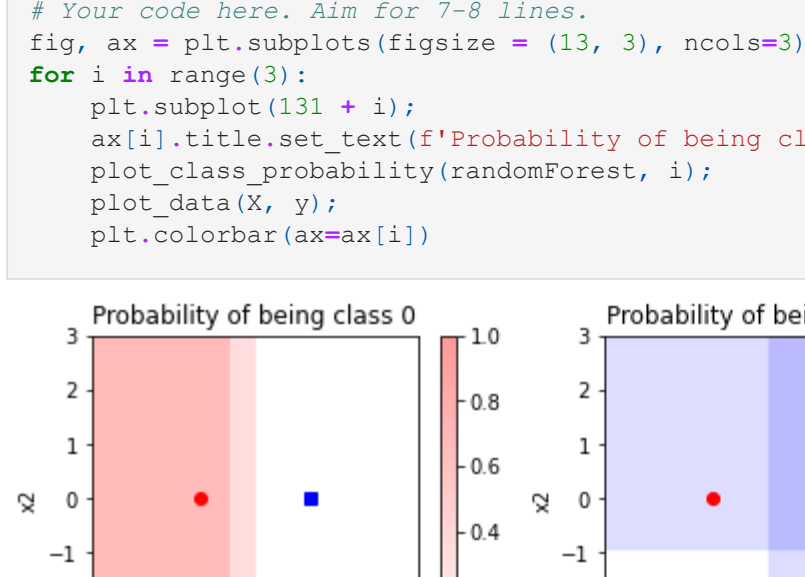
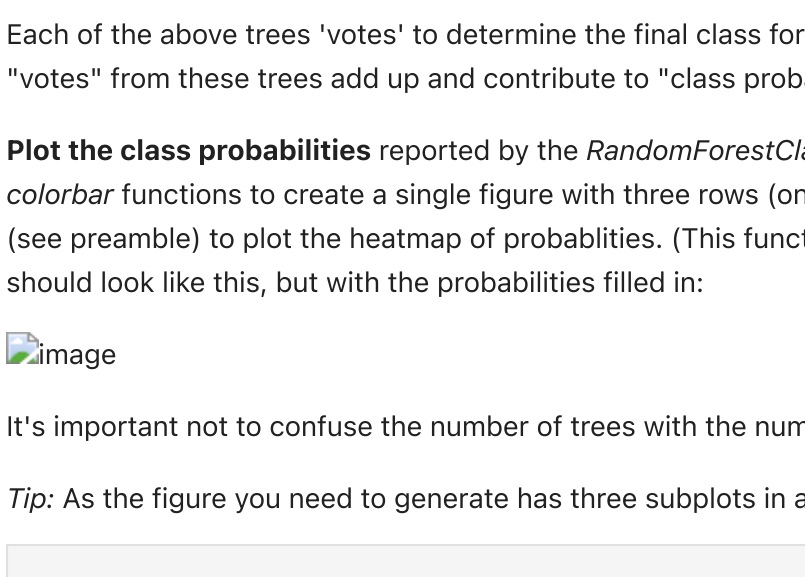
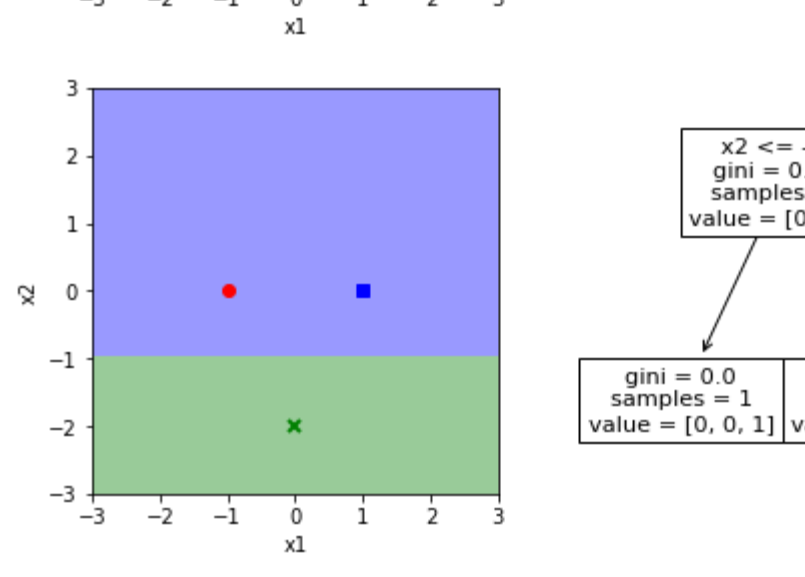
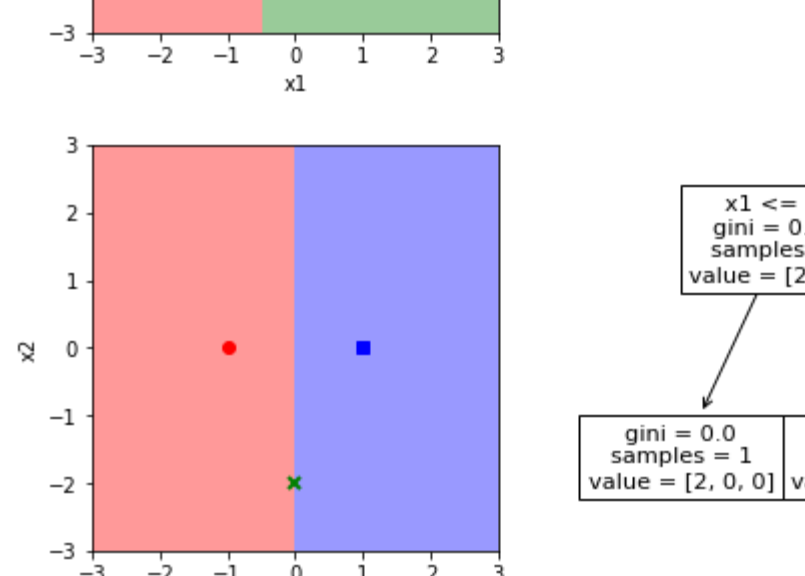
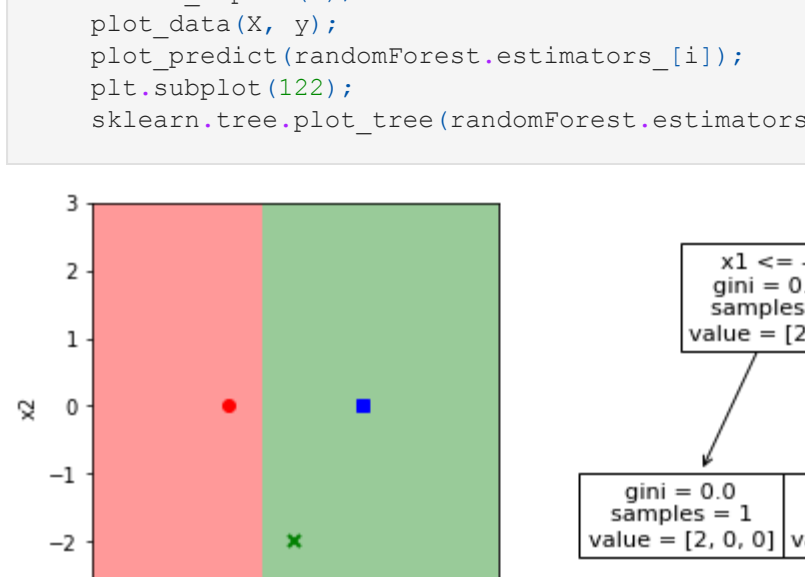
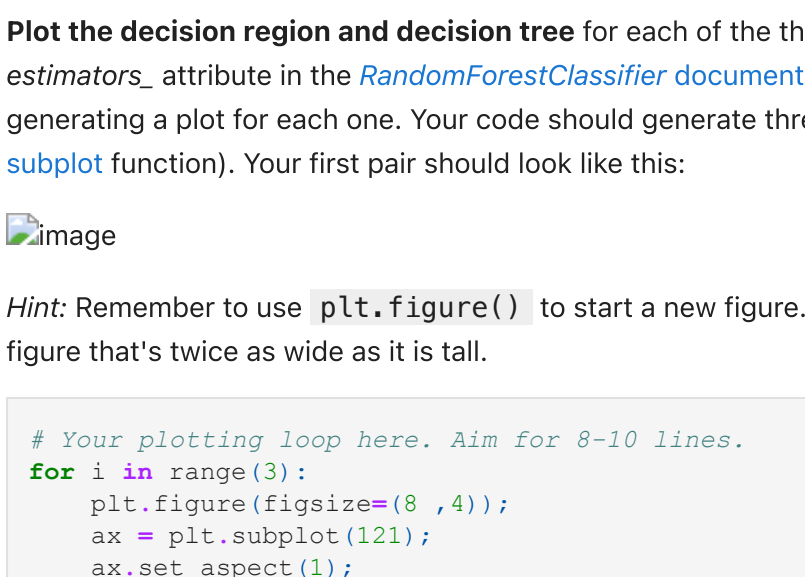
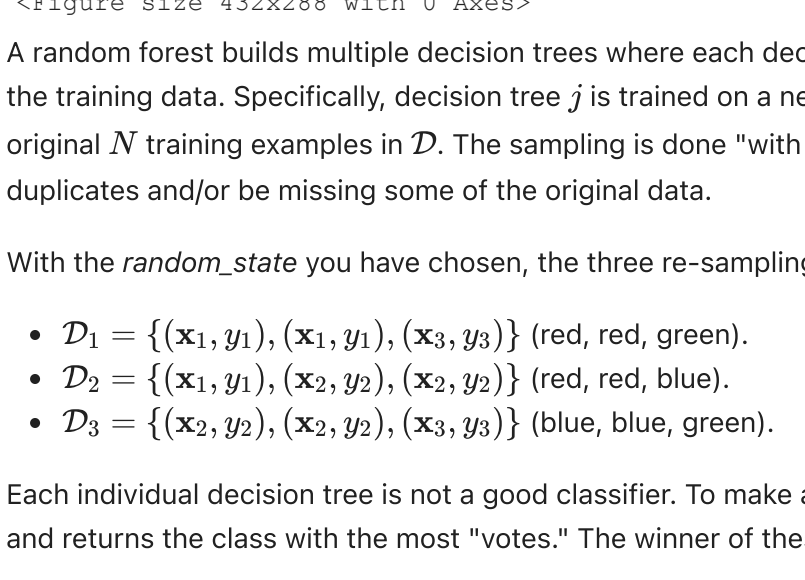
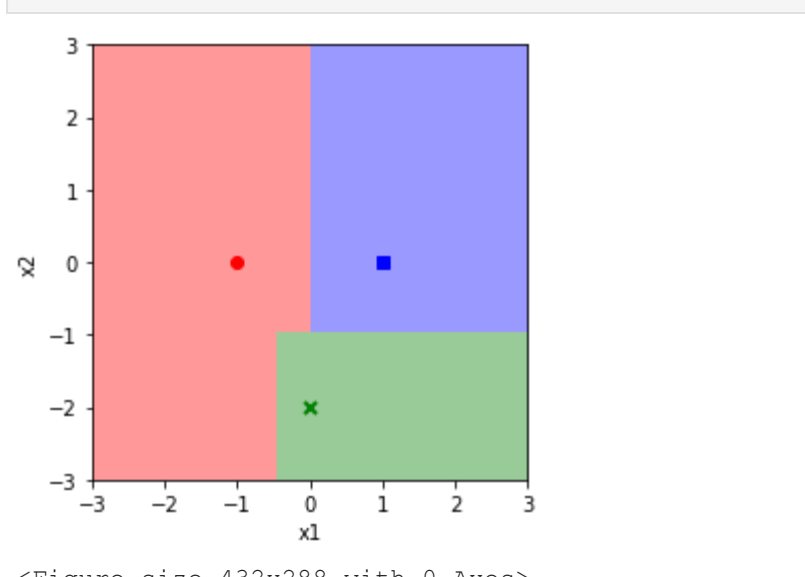
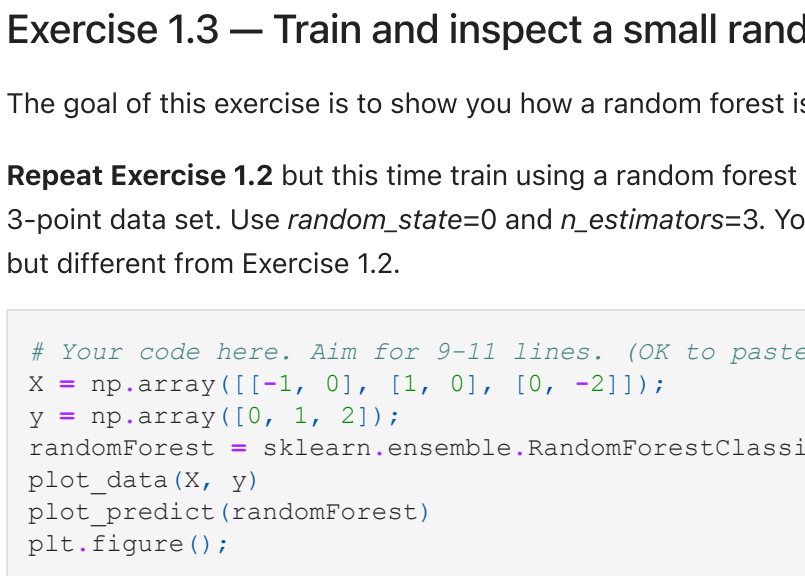
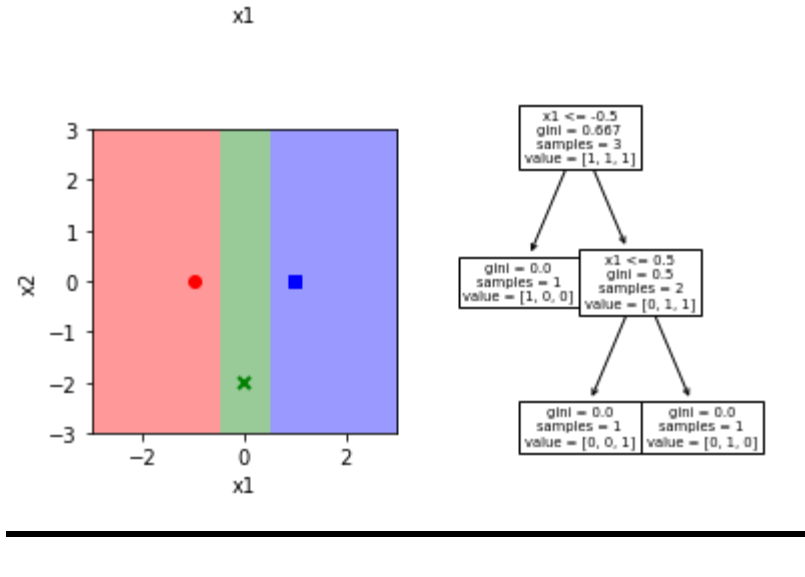
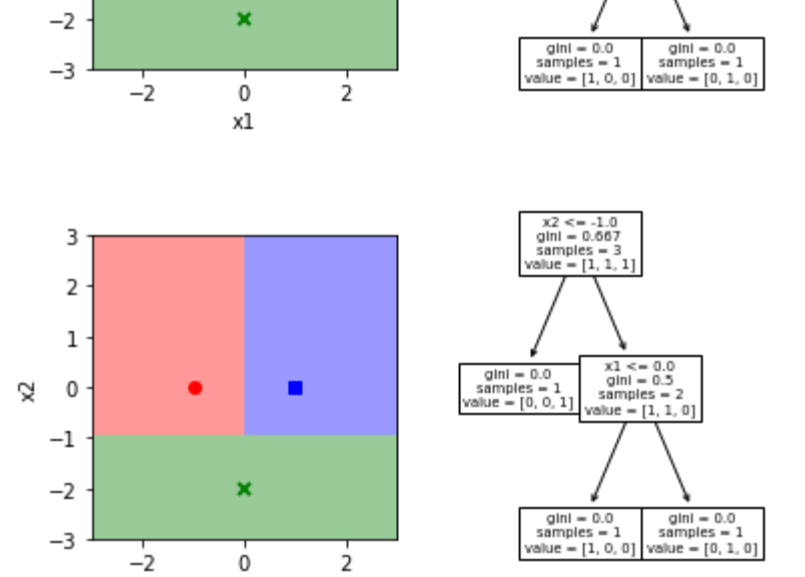
1. To the original `X`, and `x3`, add a third training point $x_3 = (0, -2)$ with class label $y_3 = 2$.
2. Print the **feature importances** attribute of your trained decision tree. Intuitively, the feature importance is 0.0 if a feature is not used at all, or 1.0 if it is the only feature needed to make decisions.

This time you should see a binary decision tree of height 2, where the first split is done by thresholding the second feature (`X[1]`) and the second split is done by thresholding the first feature (`X[0]`). If a node has value `[-1,0]`, it means that node's region (say, the red-blue region) contains exactly one training point from class 0, one training point from class 1, and zero training points from class 2.

1. Try incrementing `random_state` from 0 up to 9. How many distinct decision trees did you observe? Do the 'feature importances' make intuitive sense, given the trees that you observed?

In [44]:

```
# Your code here. Aim for 9-11 lines.
X = np.array([[ -1, 0], [ 1, 0], [0, -2]]);
y = np.array([ 0, 1, 2]);
for i in range(10):
    decisionTree = sklearn.tree.DecisionTreeClassifier(random_state=i).fit(X, y)
    print(decisionTree.feature_importances_)
    plt.figure()
    ax = plt.subplot(121)
    ax.set_aspect(1)
    plot_data(X, y)
    plot_predict(decisionTree)
    plt.subplot(122)
    sklearn.tree.plot_tree(decisionTree, feature_names=['x1', 'x2'])
```



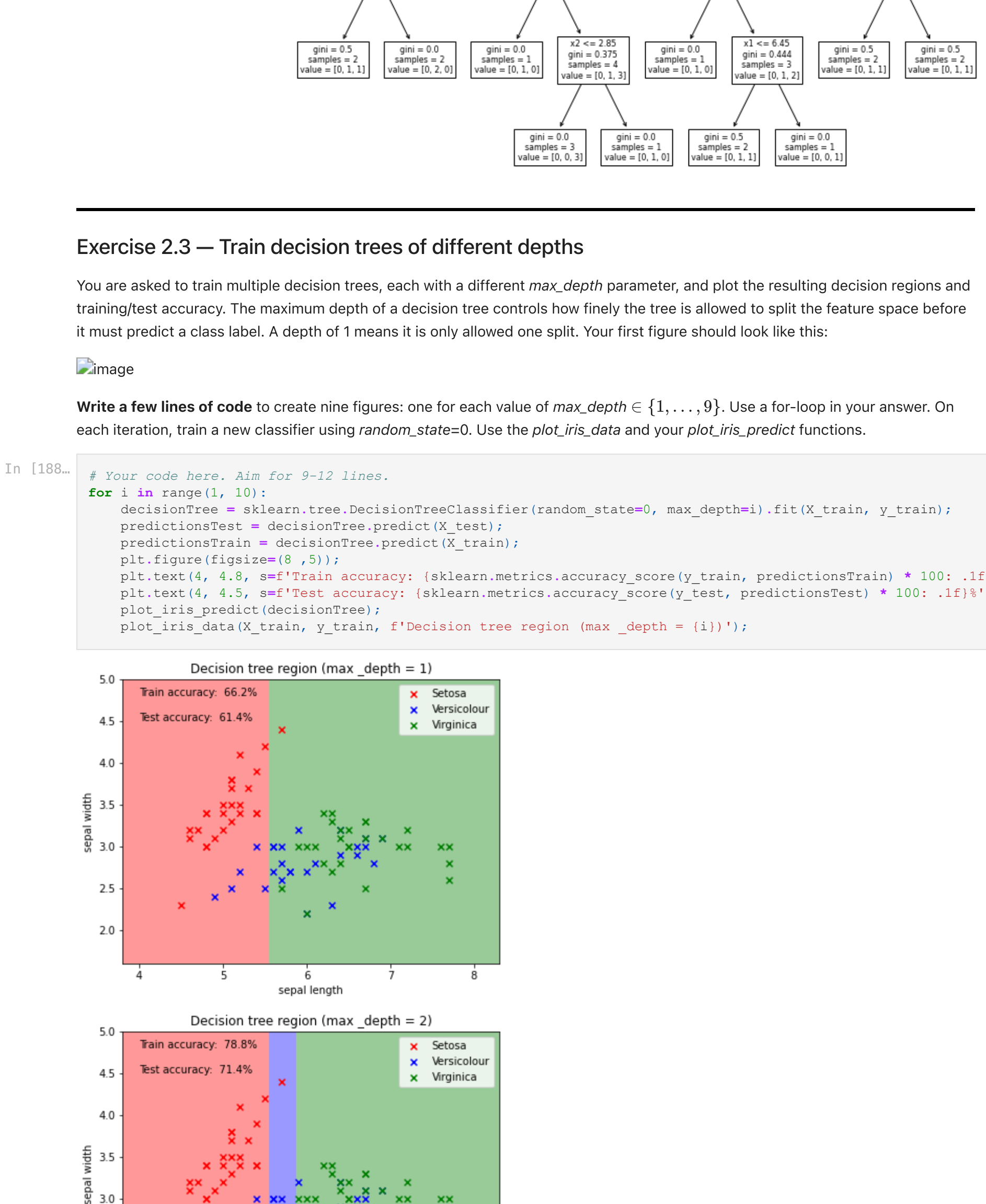

```
179...
def plot_iris_prediction(model):
    # Your prediction plotting code here. Aim for 8-12 lines.
    iris_extents = (3.5, 5.3, 1.6, 5)
    xlims, ylims, x2min, x2max = iris_extents
    x1, x2 = np.meshgrid(np.linspace(x1min, x1max, 100), np.linspace(x2min, x2max, 100))
    X = np.column_stack([x1.ravel(), x2.ravel()])
    y = model.predict(X).reshape(x1.shape)
    cmap = matplotlib.colors.ListedColormap(['r', 'b', 'g'])
    plt.imshow(y, extent=iris_extents, origin='lower', alpha=0.4, vmin=0, vmax=2, cmap=cmap, interpolation='nearest')
    plt.xlim(x1min, x1max)
    plt.ylim(x2min, x2max)
    plt.gca().set_aspect('equal')

# Your training, plotting, and evaluation code here. Aim for 7-9 lines.
decisionTree = sklearn.tree.DecisionTreeClassifier(random_state=0).fit(X_train, y_train)
predictionTest = decisionTree.predict(X_test)
predictionTrain = decisionTree.predict(X_train)
plt.figure(figsize=(8, 5))
plt.text(4, 4.8, s=f'Train accuracy: {sklearn.metrics.accuracy_score(y_train, predictionTrain) * 100:.1f}%')
plt.text(4, 4.5, s=f'Test accuracy: {sklearn.metrics.accuracy_score(y_test, predictionTest) * 100:.1f}%')
plot_iris_prediction(decisionTree)
plot_iris_data(X_train, y_train, 'Decision tree region')
```

Plot the decision tree using the `plot` tree function. You'll need to use the `figure` function with `figsize=(16,16)` in order to make the figure large enough to see all the details. Remember the hint about ending a line with a semicolon (;).

Ask yourself: is this tree reasonably "interpretable" in your view?

```
181...
# Your code here. Aim for 2 lines.
plt.figure(figsize=(16,16))
sklearn.tree.plot_tree(decisionTree, feature_names=['x1', 'x2'])
```



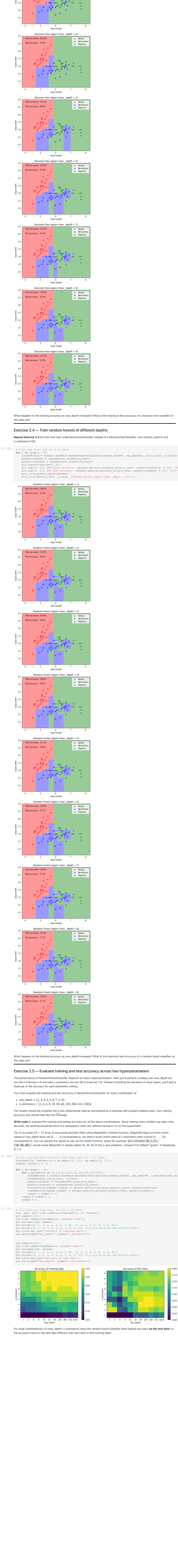
Exercise 2.3 — Train decision trees of different depths

You are asked to train multiple decision trees, each with a different `max_depth` parameter, and plot the resulting decision regions and training/test accuracy. The maximum depth of a decision tree controls how finely the tree is allowed to split the feature space before it must predict a class label. A depth of 1 means it is only allowed one split. Your first figure should look like this:

image

Write a few lines of code to create nine figures: one for each value of `max_depth` in {1,...,9}. Use a for-loop in your answer. On each iteration, train a new classifier using `random_state=0`. Use the `plot_iris_data` and your `plot_iris_prediction` functions.

```
188...
# Your code here. Aim for 9-12 lines.
for i in range(1, 10):
    decisionTree = sklearn.tree.DecisionTreeClassifier(random_state=0, max_depth=i).fit(X_train, y_train)
    predictionTest = decisionTree.predict(X_test)
    predictionTrain = decisionTree.predict(X_train)
    plt.figure(figsize=(8, 5))
    plt.text(4, 4.8, s=f'Train accuracy: {sklearn.metrics.accuracy_score(y_train, predictionTrain) * 100:.1f}%')
    plt.text(4, 4.5, s=f'Test accuracy: {sklearn.metrics.accuracy_score(y_test, predictionTest) * 100:.1f}%')
    plot_iris_prediction(decisionTree)
    plot_iris_data(X_train, y_train, f'Decision tree region (max_depth = {i})')
```



What happens to the training accuracy as `max_depth` increases? What is the maximum test accuracy of a decision tree classifier on this data set?

Exercise 2.4 — Train random forests of different depths

Repeat Exercise 2.3 but this time train a `RandomForestClassifier` instead of a `DecisionTreeClassifier`. Use `random_state=0` and `n_estimators=100`.

```
189...
# Your code here. Aim for 9-12 lines.
for i in range(1, 10):
    randomForest = sklearn.ensemble.RandomForestClassifier(random_state=0, max_depth=i).fit(X_train, y_train)
    predictionTest = randomForest.predict(X_test)
    predictionTrain = randomForest.predict(X_train)
    plt.figure(figsize=(8, 5))
    plt.text(4, 4.8, s=f'Train accuracy: {sklearn.metrics.accuracy_score(y_train, predictionTrain) * 100:.1f}%')
    plt.text(4, 4.5, s=f'Test accuracy: {sklearn.metrics.accuracy_score(y_test, predictionTest) * 100:.1f}%')
    plot_iris_prediction(randomForest)
    plot_iris_data(X_train, y_train, f'Random Forest region (max_depth = {i})')
```



What happens to the training accuracy as `max_depth` increases? What is the maximum test accuracy of a random forest classifier on this data set?

Exercise 2.5 — Evaluate training and test accuracy across two hyperparameters

The performance of `RandomForestClassifier` depends on many hyperparameters. Here you'll perform a sweep over `max_depth` (as you did in Exercise 2.4) and also `n_estimators` (as you did in Exercise 1.4). Instead of plotting the decisions in input space, you'll plot a heatmap of the accuracy for each parameter setting.

You must evaluate the training and test accuracy of `RandomForestClassifier` for every combination of

- `max_depth` in {1, 2, 3, 4, 5, 6, 7, 8, 9}
- `n_estimators` in {1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024}

The results should be compiled into a two-dimensional ndarray and plotted as a heatmap with properly labelled axes. Your training accuracy plot should look like this: image

Write code to evaluate the training and testing accuracy for all the above combinations. Since training many models may take a few seconds, the plotting should be done in a subsequent code cell, without having to re-run the experiment.

Tip: you build a 9×11 array of accuracies and plot them using Matplotlib's `imshow` function. Matplotlib does not know which values of `max_depth` each row 0,...,8 corresponds to, nor does it know which value of `n_estimators` each column 0,...,10 corresponds to. You can specify the values to use via the `xticks` function, where for example `plt.xticks([0,1,2], [10,20,40])` would cause Matplotlib to display labels 10, 20, 40 at the x-axis positions, instead of its default "guess" of 0, 1, 2.

```
190...
# Your training and evaluation code here. Aim for 8-12 lines.
trainResults, testResults = np.empty((9, 11), np.empty((9, 11)))
indexX = 0;
indexY = 0;
for i in range(1, 10):
    for n_estimators in [1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024]:
        randomForest = sklearn.ensemble.RandomForestClassifier(random_state=0, max_depth=i, n_estimators=n_estimators)
        predictionTest = randomForest.predict(X_test)
        predictionTrain = randomForest.predict(X_train)
        trainResults[indexX][indexY] = sklearn.metrics.accuracy_score(y_train, predictionTrain)
        testResults[indexX][indexY] = sklearn.metrics.accuracy_score(y_test, predictionTest)
        indexX = indexX + 1
        indexY = 0;
```

```
192...
# Your plotting code here. Aim for 8-10 lines.
fig, (ax1, ax2) = plt.subplots(figsize=(15, 5), ncols=2);
plt.subplot(121);
lab = ax1.imshow(trainResults, origin='lower');
plt.colorbar(lab, ax=ax1);
plt.yticks([0, 1, 2, 3, 4, 5, 6, 7, 8], [1, 2, 3, 4, 5, 6, 7, 8, 9]);
plt.xticks([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10], [1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024]);
ax1.title.set_text('Accuracy of training data');
ax1.set_xlabel('max_depth', ylabel='n_estimators');
```



For what combination(s) of (`max_depth`, `n_estimators`) does the random forest classifier have highest accuracy on the test data? Is the accuracy trend on the test data different than the trend on the training data?