# ML-Assignment 1 - Mathieu Rio, Remi Maigrot

January 29, 2025

# 1 Machine Learning - Assignment 1

## 1.1 Naive Bayes Learning algorithm, Cross-validation, and ROC-Curves

The aim of the assignment is to implement:

- Naive Bayes learning algorithm for binary classification tasks
- Visualization to plot a ROC-curve
- A cross-validation test
- Visualization of the average ROC-curve of a cross-validation test

Follow the instructions and implement what is missing to complete the assignment. Some functions have been started to help you a little bit with the inputs or outputs of the function.

**Note:** You might need to go back and forth during your implementation of the code. The structure is set up to make implementation easier, but how you return values from the different functions might vary, and you might find yourself going back and change something to make it easier later on.

## 1.2 Assignment preparations

We help you out with importing the libraries and reading the data.

Look at the output to get an idea of how the data is structured.

```python
[118]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from math import e, pi, sqrt
```

```python
[119]: class Data_set:
    def __init__(self):
        self.data = None
        self.features = None
        self.labels = None
        self.means = 0.0
        self.prior = 0.0
        self.std_devs = 0.0
        self.gaussian_probability_density = 0.0
```

```python
    def fix_data_structure(self):
        if all(isinstance(row, np.ndarray) for row in self.features):
            try:
                self.features = np.array(self.features, dtype=float)
            except ValueError as e:
                print("Erreur lors de la conversion :", e)
        else:
            print("Les éléments ne sont pas tous des tableaux numpy.")

    def class_split(self):
        self.features = self.data[:, :-1]
        self.labels = self.data[:, -1]
        self.fix_data_structure()

    def display_features_and_labels(self):
        print("Features set:")
        print(self.features)
        print("Labels set:")
        print(self.labels)
```

```python
[120]: # creating a class to make the code cleaner
class Flower:
    def __init__(self):
        self.data = None
        self.train = Data_set()
        self.test = Data_set()
        self.test_size = 0.2

    def train_test_split(self):
        np.random.shuffle(data)

        split_index = int(len(self.data) * (1 - self.test_size))

        self.train.data = self.data[:split_index]
        self.test.data = self.data[split_index:]

    def class_split_automation(self):
        self.train.class_split()
        self.test.class_split()
        self.train.display_features_and_labels()
        self.test.display_features_and_labels()

    def display_training(self):
        print("Train set:")
        print(self.train.data[:3])
        print("Test set:")
        print(self.test.data[:3])
```

```
[121]: data = pd.read_csv("iris.csv").to_numpy()

       mapped, index, unique_arr = np.unique(data[:, -1], return_index=True,␣
        ↪return_inverse=True)
       data[:, -1] = unique_arr

       iris_setosa = Flower()
       iris_versicolor = Flower()
       iris_virginica = Flower()

       iris_setosa.data, iris_versicolor.data, iris_virginica.data = np.split(data,␣
        ↪index[1:])

       print(f"Full data array (features and labels):\n{iris_setosa.data[:3]}\n")
       print("###############\n")
       print(f"Train features (first 4 columns):\n{iris_setosa.data[:3, :-1]}\n")
       print(f"Labels (last column):\n{iris_setosa.data[:3, -1:]}\n")
       print(f"Names of labels:\n{[[numb, name] for numb, name in enumerate(mapped)]}")
```

```
Full data array (features and labels):
[[5.1 3.5 1.4 0.2 0]
 [4.9 3.0 1.4 0.2 0]
 [4.7 3.2 1.3 0.2 0]]

###############

Train features (first 4 columns):
[[5.1 3.5 1.4 0.2]
 [4.9 3.0 1.4 0.2]
 [4.7 3.2 1.3 0.2]]

Labels (last column):
[[0]
 [0]
 [0]]

Names of labels:
[[0, 'Iris-setosa'], [1, 'Iris-versicolor'], [2, 'Iris-virginica']]
```

```
[122]: # Example print of the 3 first datapoints (similar as above):
       iris_setosa.data[:3]
```

```
[122]: array([[5.1, 3.5, 1.4, 0.2, 0],
              [4.9, 3.0, 1.4, 0.2, 0],
              [4.7, 3.2, 1.3, 0.2, 0]], dtype=object)
```

## 1.3 Data handling functions

As a start, we are going to implement some basic data handling functions to use in the future.

### 1.3.1 1) Split class into a train and test set

First, we need to be able to split the class into a train and test set.

```
[123]: # TODO: Test the train_test_split function
iris_setosa.train_test_split()
iris_versicolor.train_test_split()
iris_virginica.train_test_split()

# TODO: Print the output
iris_setosa.display_training()
iris_versicolor.display_training()
iris_virginica.display_training()
```

```
Train set:
[[6.7 3.1 4.4 1.4 1]
 [5.0 3.0 1.6 0.2 0]
 [5.5 2.3 4.0 1.3 1]]
Test set:
[[5.1 3.5 1.4 0.2 0]
 [5.1 3.8 1.5 0.3 0]
 [6.8 3.2 5.9 2.3 2]]
Train set:
[[7.7 2.6 6.9 2.3 2]
 [5.7 2.5 5.0 2.0 2]
 [4.9 2.4 3.3 1.0 1]]
Test set:
[[7.1 3.0 5.9 2.1 2]
 [6.2 2.9 4.3 1.3 1]
 [6.7 3.0 5.2 2.3 2]]
Train set:
[[6.1 2.8 4.7 1.2 1]
 [6.7 3.1 5.6 2.4 2]
 [6.4 2.9 4.3 1.3 1]]
Test set:
[[5.0 2.0 3.5 1.0 1]
 [5.0 3.4 1.5 0.2 0]
 [6.5 3.0 5.8 2.2 2]]
```

### 1.3.2 2) Split data into features and labels

The data as shown above is not always the optimal shape. To help us keep track of things, we can split the data into its features and labels seperately.

Each class is 4 features and 1 label in the same array:

- **[feature 1, feature 2, feature 3, feature 4, label]**

It would help us later to have the features and labels in seperate arrays in the form:

- **[feature 1, feature 2, feature 3, feature 4]** and **[label]**

Here you are going to implement this functionallity.

We should first test the "**class_split**" function on one of the classes above (iris_setosa, etc...) to make sure it works properly.

```
[124]:  # TODO: Test the class splitting function
        # iris_setosa.train.class_split()
        # iris_setosa.test.class_split()

        # iris_versicolor.train.class_split()
        # iris_versicolor.test.class_split()

        # iris_virginica.train.class_split()
        # iris_virginica.test.class_split()

        # TODO: Print the output
        # iris_setosa.train.display_features_and_labels()
        # iris_setosa.test.display_features_and_labels()

        # iris_versicolor.train.display_features_and_labels()
        # iris_versicolor.test.display_features_and_labels()

        # iris_virginica.train.display_features_and_labels()
        # iris_virginica.test.display_features_and_labels()

        # Or

        iris_setosa.class_split_automation()
        iris_versicolor.class_split_automation()
        iris_virginica.class_split_automation()
```

```
Features set:
[[6.7 3.1 4.4 1.4]
 [5.  3.  1.6 0.2]
 [5.5 2.3 4.  1.3]
 [6.4 2.8 5.6 2.1]
 [7.2 3.6 6.1 2.5]
 [6.7 3.  5.  1.7]
 [6.4 2.8 5.6 2.2]
 [4.3 3.  1.1 0.1]
 [4.6 3.1 1.5 0.2]
 [5.8 2.7 4.1 1. ]
 [5.1 3.8 1.9 0.4]
 [4.7 3.2 1.3 0.2]
```

```
 [5.1 3.5 1.4 0.3]
 [6.  2.2 5.  1.5]
 [7.7 3.8 6.7 2.2]
 [5.6 3.  4.1 1.3]
 [5.2 3.4 1.4 0.2]
 [5.  3.4 1.6 0.4]
 [6.9 3.2 5.7 2.3]
 [5.9 3.2 4.8 1.8]
 [7.6 3.  6.6 2.1]
 [6.8 3.  5.5 2.1]
 [4.9 3.1 1.5 0.1]
 [4.7 3.2 1.6 0.2]
 [7.2 3.2 6.  1.8]
 [5.4 3.9 1.3 0.4]
 [5.4 3.  4.5 1.5]
 [5.8 4.  1.2 0.2]
 [6.7 2.5 5.8 1.8]
 [6.3 2.7 4.9 1.8]
 [4.4 2.9 1.4 0.2]
 [7.3 2.9 6.3 1.8]
 [5.4 3.9 1.7 0.4]
 [6.9 3.1 5.4 2.1]
 [7.7 3.  6.1 2.3]
 [6.7 3.1 4.7 1.5]
 [5.7 2.8 4.5 1.3]
 [4.8 3.4 1.9 0.2]
 [5.5 2.4 3.7 1. ]
 [6.3 3.3 4.7 1.6]]
Labels set:
[1 0 1 2 2 1 2 0 0 1 0 0 0 2 2 1 0 0 2 1 2 2 0 0 2 0 1 0 2 2 0 2 0 2 2 1 1
 0 1 1]
Features set:
[[5.1 3.5 1.4 0.2]
 [5.1 3.8 1.5 0.3]
 [6.8 3.2 5.9 2.3]
 [5.8 2.6 4.  1.2]
 [5.  3.2 1.2 0.2]
 [6.9 3.1 4.9 1.5]
 [5.  3.3 1.4 0.2]
 [4.4 3.  1.3 0.2]
 [4.9 3.1 1.5 0.1]
 [6.3 2.5 5.  1.9]]
Labels set:
[0 0 2 1 0 1 0 0 0 2]
Features set:
[[7.7 2.6 6.9 2.3]
 [5.7 2.5 5.  2. ]
 [4.9 2.4 3.3 1. ]
```

```
 [5.7 2.9 4.2 1.3]
 [7.2 3.  5.8 1.6]
 [6.4 3.2 4.5 1.5]
 [5.1 3.3 1.7 0.5]
 [6.3 2.9 5.6 1.8]
 [6.5 3.2 5.1 2. ]
 [6.4 3.1 5.5 1.8]
 [6.7 3.3 5.7 2.5]
 [6.3 3.4 5.6 2.4]
 [5.8 2.7 5.1 1.9]
 [6.1 2.9 4.7 1.4]
 [4.6 3.2 1.4 0.2]
 [5.7 3.  4.2 1.2]
 [5.7 2.8 4.1 1.3]
 [5.6 2.8 4.9 2. ]
 [6.  2.9 4.5 1.5]
 [5.8 2.7 5.1 1.9]
 [6.4 2.7 5.3 1.9]
 [6.1 2.6 5.6 1.4]
 [4.9 3.  1.4 0.2]
 [4.9 2.5 4.5 1.7]
 [4.4 3.2 1.3 0.2]
 [5.1 2.5 3.  1.1]
 [6.7 3.3 5.7 2.1]
 [5.2 3.5 1.5 0.2]
 [5.8 2.8 5.1 2.4]
 [6.5 3.  5.2 2. ]
 [5.5 3.5 1.3 0.2]
 [6.2 3.4 5.4 2.3]
 [7.4 2.8 6.1 1.9]
 [5.8 2.7 3.9 1.2]
 [6.3 3.3 6.  2.5]
 [5.1 3.8 1.6 0.2]
 [5.7 3.8 1.7 0.3]
 [5.7 4.4 1.5 0.4]
 [5.5 2.6 4.4 1.2]
 [5.4 3.4 1.7 0.2]]
Labels set:
[2 2 1 1 2 1 0 2 2 2 2 2 2 1 0 1 1 2 1 2 2 2 0 2 0 1 2 0 2 2 0 2 2 1 2 0 0
 0 1 0]
Features set:
[[7.1 3.  5.9 2.1]
 [6.2 2.9 4.3 1.3]
 [6.7 3.  5.2 2.3]
 [5.9 3.  5.1 1.8]
 [4.9 3.1 1.5 0.1]
 [7.  3.2 4.7 1.4]
 [6.1 3.  4.9 1.8]
```

```
 [5.  2.3 3.3 1. ]
 [5.6 2.7 4.2 1.3]
 [4.8 3.  1.4 0.3]]
Labels set:
[2 1 2 2 0 1 2 1 1 0]
Features set:
[[6.1 2.8 4.7 1.2]
 [6.7 3.1 5.6 2.4]
 [6.4 2.9 4.3 1.3]
 [5.6 2.9 3.6 1.3]
 [6.1 2.8 4.  1.3]
 [5.  3.5 1.3 0.3]
 [4.6 3.4 1.4 0.3]
 [5.6 2.5 3.9 1.1]
 [6.5 2.8 4.6 1.5]
 [5.1 3.7 1.5 0.4]
 [7.9 3.8 6.4 2. ]
 [5.5 2.4 3.8 1.1]
 [6.3 2.8 5.1 1.5]
 [6.  3.  4.8 1.8]
 [5.2 2.7 3.9 1.4]
 [6.3 2.3 4.4 1.3]
 [4.6 3.6 1.  0.2]
 [5.4 3.7 1.5 0.2]
 [6.4 3.2 5.3 2.3]
 [5.  3.5 1.6 0.6]
 [5.7 2.6 3.5 1. ]
 [6.9 3.1 5.1 2.3]
 [5.3 3.7 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.1 3.4 1.5 0.2]
 [6.5 3.  5.5 1.8]
 [6.6 3.  4.4 1.4]
 [5.5 4.2 1.4 0.2]
 [5.6 3.  4.5 1.5]
 [6.  3.4 4.5 1.6]
 [5.5 2.5 4.  1.3]
 [5.4 3.4 1.5 0.4]
 [4.8 3.1 1.6 0.2]
 [6.  2.7 5.1 1.6]
 [4.5 2.3 1.3 0.3]
 [6.6 2.9 4.6 1.3]
 [7.7 2.8 6.7 2. ]
 [5.2 4.1 1.5 0.1]
 [4.8 3.  1.4 0.1]
 [4.8 3.4 1.6 0.2]]
Labels set:
[1 2 1 1 1 0 0 1 1 0 2 1 2 2 1 1 0 0 2 0 1 2 0 0 0 2 1 0 1 1 1 0 0 1 0 1 2
```

8

```
 0 0 0]
Features set:
[[5.   2.   3.5 1. ]
 [5.   3.4 1.5 0.2]
 [6.5 3.   5.8 2.2]
 [6.8 2.8 4.8 1.4]
 [5.9 3.   4.2 1.5]
 [6.2 2.2 4.5 1.5]
 [6.   2.2 4.   1. ]
 [6.1 3.   4.6 1.4]
 [6.3 2.5 4.9 1.5]
 [6.2 2.8 4.8 1.8]]
Labels set:
[1 0 2 1 1 1 1 1 1 2]
```

We should also try to **1)** first split a class into a train and test set, **2)** split each of these two into features and abels. In total there should be 4 arrays (2 feature and 2 label arrays). Think a bit before going to the next task, what can easily go wrong in the above code?

## 1.4 Naive Bayes learning algorithm

When implementing the Navie Bayes learning algorithm, we can break it down into a few components.

We will implement these components one at a time.

### 1.4.1 3) Calculate feature statistics

First, we need to implement a function that returns feature statistics (means, standard deviation, priors) for a given set of feature data for a single class. This is the equivalent of "training" the naive bayes model.

**Note 1:** Each feature gets its own mean and standard deviation!

**Note 2:** The way you structure the functions (what is returned) shapes the remainder of the assignment.

```python
[125]: def calculate_feature_statistics(data, total):
           data = np.array(data, dtype=float)

           means = np.mean(data, axis=0)
           # means = [round(val, 3) for val in means]

           std_devs = np.std(data, axis=0, ddof=0)
           # std_devs = [round(val, 3) for val in std_devs]

           prior = len(data) / len(total)
           # prior = [round(val, 3) for val in prior]

           return means, std_devs, prior
```

To make sure the function works, we should test it before proceding.

```python
[126]: # TODO: Make sure to use our previous class splitting function.
       # print(type(iris_setosa.train.features), iris_setosa.train.features)

       # TODO: Test the function here for one of the dataset classes.
       iris_setosa.train.means, iris_setosa.train.std_devs, iris_setosa.train.prior =␣
        ↪calculate_feature_statistics(iris_setosa.train.features, iris_setosa.data)
       iris_versicolor.train.means, iris_versicolor.train.std_devs, iris_versicolor.
        ↪train.prior = calculate_feature_statistics(iris_versicolor.train.features,␣
        ↪iris_versicolor.data)
       iris_virginica.train.means, iris_virginica.train.std_devs, iris_virginica.train.
        ↪prior = calculate_feature_statistics(iris_virginica.train.features,␣
        ↪iris_virginica.data)

       # TODO: Print the output from the feature statistic function.
       # print("feature", iris_setosa.train.features)
       print("means", iris_setosa.train.means)
       print("std_devs", iris_setosa.train.std_devs)
       print("prior", iris_setosa.train.prior)
```

```
means [5.9325 3.1125 3.805  1.1925]
std_devs [0.96264934 0.41484184 1.91245261 0.80261681]
prior 0.8
```

### 1.4.2  4) Gaussian probability density function (Gaussian PDF)

Now we need to implement the gaussian probability density function to use for a single datapoint.

**Note:** Look at the imports in the first cell at the top, it has some math numbers for easy use here.

```python
[127]: def gaussian_probability_density_function(x, mean, stdev):
           exponent = np.exp(-((x - mean) ** 2) / (2 * (stdev ** 2)))
           coefficient = 1 / (np.sqrt(2 * np.pi * (stdev ** 2)))
           return coefficient * exponent
```

```python
[128]: def for_each_gaussian(item):
           item.train.gaussian_probability_density = [
               [
                   gaussian_probability_density_function(x, mean, stdev)
                   for x, mean, stdev in zip(row, item.train.means, item.train.
        ↪std_devs)
               ]
               for row in item.train.features
           ]
```

### 1.4.3 5) Testing Gaussian PDF

We should test it to make sure it works. Train it, using the "calculate_feature_statistics" function, on one of the dataset classes. Then, take one datapoint from the same class and use naive bayes gaussian to make a prediction.

[129]:
```python
# TODO: Implement the code below to test the␣
 ↪"gaussian_probability_density_function" function for one of the classes.
# TODO: Test with one datapoint from the learned class.

for_each_gaussian(iris_setosa)
for_each_gaussian(iris_versicolor)
for_each_gaussian(iris_virginica)

#print(iris_setosa.train.gaussian_probability_density)

# TODO: Print the probability density
for n in iris_setosa.train.gaussian_probability_density:
        print(n)
```

```
[np.float64(0.3015862352656024), np.float64(0.9612367382441159),
np.float64(0.19874702333108324), np.float64(0.4807156381694144)]
[np.float64(0.2592289112248091), np.float64(0.9269534144089494),
np.float64(0.10731416151553415), np.float64(0.2313957423797314)]
[np.float64(0.37463646594007766), np.float64(0.1412683473866337),
np.float64(0.20752088352629838), np.float64(0.49261358947966627)]
[np.float64(0.36832298338079217), np.float64(0.7241111434083813),
np.float64(0.1342842752239309), np.float64(0.2622976432133863)]
[np.float64(0.1741751234952665), np.float64(0.4821184614791616),
np.float64(0.10153408512159703), np.float64(0.13186762713571298)]
[np.float64(0.3015862352656024), np.float64(0.9269534144089494),
np.float64(0.1716075482309766), np.float64(0.40699004582329756)]
[np.float64(0.36832298338079217), np.float64(0.7241111434083813),
np.float64(0.1342842752239309), np.float64(0.22606995688004783)]
[np.float64(0.09839055839869965), np.float64(0.9269534144089494),
np.float64(0.07671878324957038), np.float64(0.1968218924340225)]
[np.float64(0.1589956533773653), np.float64(0.9612367382441159),
np.float64(0.10089759229016843), np.float64(0.2313957423797314)]
[np.float64(0.4105141068117921), np.float64(0.5865769674241764),
np.float64(0.20613542747756042), np.float64(0.4829595440764669)]
[np.float64(0.2851292334366234), np.float64(0.24357440672122965),
np.float64(0.12701682165688594), np.float64(0.30527712310590566)]
[np.float64(0.1825952765984702), np.float64(0.9405175189493056),
np.float64(0.08846390105810852), np.float64(0.2313957423797314)]
[np.float64(0.2851292334366234), np.float64(0.6216714600761556),
np.float64(0.09460566720451268), np.float64(0.2678524662451183)]
[np.float64(0.4134036501282108), np.float64(0.0855825098965548),
np.float64(0.1716075482309766), np.float64(0.46187923475148857)]
```

11

```
[np.float64(0.07680657249284024), np.float64(0.24357440672122965),
np.float64(0.06633323350742046), np.float64(0.22606995688004783)]
[np.float64(0.3904234438889966), np.float64(0.9269534144089494),
np.float64(0.20613542747756042), np.float64(0.49261358947966627)]
[np.float64(0.3102512665079876), np.float64(0.7563663060461174),
np.float64(0.09460566720451268), np.float64(0.2313957423797314)]
[np.float64(0.2592289112248091), np.float64(0.7563663060461174),
np.float64(0.10731416151553415), np.float64(0.30527712310590566)]
[np.float64(0.2500926007228404), np.float64(0.9405175189493056),
np.float64(0.1276783706450599), np.float64(0.19184462667657717)]
[np.float64(0.4141850732442614), np.float64(0.9405175189493056),
np.float64(0.1821968230133562), np.float64(0.373249255330593)]
[np.float64(0.09244616967779351), np.float64(0.9269534144089494),
np.float64(0.07169900966824823), np.float64(0.2622976432133863)]
[np.float64(0.2761209929395077), np.float64(0.9269534144089494),
np.float64(0.14084634162058496), np.float64(0.2622976432133863)]
[np.float64(0.2331517239253164), np.float64(0.9612367382441159),
np.float64(0.10089759229016843), np.float64(0.1968218924340225)]
[np.float64(0.1825952765984702), np.float64(0.9405175189493056),
np.float64(0.10731416151553415), np.float64(0.2313957423797314)]
[np.float64(0.1741751234952665), np.float64(0.9405175189493056),
np.float64(0.10796160991803662), np.float64(0.373249255330593)]
[np.float64(0.35562944645979666), np.float64(0.15867799816733408),
np.float64(0.08846390105810852), np.float64(0.30527712310590566)]
[np.float64(0.35562944645979666), np.float64(0.9269534144089494),
np.float64(0.19527282032376272), np.float64(0.46187923475148857)]
[np.float64(0.4105141068117921), np.float64(0.0975362037654661),
np.float64(0.08249499578351092), np.float64(0.2313957423797314)]
[np.float64(0.3015862352656024), np.float64(0.32333802672114853),
np.float64(0.12106597056977161), np.float64(0.373249255330593)]
[np.float64(0.38529638356665563), np.float64(0.5865769674241764),
np.float64(0.17706482775110818), np.float64(0.373249255330593)]
[np.float64(0.11671245246582783), np.float64(0.8434309465388058),
np.float64(0.09460566720451268), np.float64(0.2313957423797314)]
[np.float64(0.15109203374719313), np.float64(0.8434309465388058),
np.float64(0.08907065114176908), np.float64(0.373249255330593)]
[np.float64(0.35562944645979666), np.float64(0.15867799816733408),
np.float64(0.11382714793885), np.float64(0.30527712310590566)]
[np.float64(0.2500926007228404), np.float64(0.9612367382441159),
np.float64(0.14732571801327685), np.float64(0.2622976432133863)]
[np.float64(0.07680657249284024), np.float64(0.9269534144089494),
np.float64(0.10153408512159703), np.float64(0.19184462667657717)]
[np.float64(0.3015862352656024), np.float64(0.9612367382441159),
np.float64(0.18696567522223922), np.float64(0.46187923475148857)]
[np.float64(0.4025086648588843), np.float64(0.7241111434083813),
np.float64(0.19527282032376272), np.float64(0.49261358947966627)]
[np.float64(0.20744708281699475), np.float64(0.7563663060461174),
np.float64(0.12701682165688594), np.float64(0.2313957423797314)]
```

```
    [np.float64(0.37463646594007766), np.float64(0.22002335676600374),
    np.float64(0.20828827322041524), np.float64(0.4829595440764669)]
    [np.float64(0.38529638356665563), np.float64(0.8682953871969392),
    np.float64(0.18696567522223922), np.float64(0.4369451769503971)]
```

[130]:
```python
for n in iris_versicolor.train.gaussian_probability_density:
    print(n)
```

```
[np.float64(0.022229597314218495), np.float64(0.5477078274217113),
np.float64(0.06046065281004332), np.float64(0.2569007773342541)]
[np.float64(0.5375348536877407), np.float64(0.41117733955798913),
np.float64(0.20805641059740262), np.float64(0.38189768467369134)]
[np.float64(0.22406173998598164), np.float64(0.2911368627189683),
np.float64(0.21091606155113538), np.float64(0.46067757337317483)]
[np.float64(0.5375348536877407), np.float64(0.9112448342915718),
np.float64(0.23810534192895577), np.float64(0.5230195407313964)]
[np.float64(0.1012353909951289), np.float64(0.9605192178192339),
np.float64(0.14467596495987176), np.float64(0.5075139966010033)]
[np.float64(0.42208106874816853), np.float64(0.8953883092545903),
np.float64(0.2325001490145239), np.float64(0.5216526317547565)]
[np.float64(0.31288342509729494), np.float64(0.7918552521983719),
np.float64(0.08327280126275306), np.float64(0.2630231827356142)]
[np.float64(0.4628070552230153), np.float64(0.9112448342915718),
np.float64(0.16186101321783258), np.float64(0.4558805058360182)]
[np.float64(0.3776222618737363), np.float64(0.8953883092545903),
np.float64(0.20131810361017977), np.float64(0.38189768467369134)]
[np.float64(0.42208106874816853), np.float64(0.9549152709974453),
np.float64(0.17029043937507005), np.float64(0.4558805058360182)]
[np.float64(0.28535038299063187), np.float64(0.7918552521983719),
np.float64(0.15330075552164707), np.float64(0.1807565062254383)]
[np.float64(0.4628070552230153), np.float64(0.6604926430655316),
np.float64(0.16186101321783258), np.float64(0.21737923326444122)]
[np.float64(0.550056946262832), np.float64(0.6881078683113983),
np.float64(0.20131810361017977), np.float64(0.4209082106732165)]
[np.float64(0.5252978276017184), np.float64(0.9112448342915718),
np.float64(0.22478960608145954), np.float64(0.5269119316680196)]
[np.float64(0.11758150717984477), np.float64(0.8953883092545903),
np.float64(0.06319073974654746), np.float64(0.15241623272983265)]
[np.float64(0.5375348536877407), np.float64(0.9605192178192339),
np.float64(0.23810534192895577), np.float64(0.510177203090685)]
[np.float64(0.5375348536877407), np.float64(0.8153646920821799),
np.float64(0.23829661722185952), np.float64(0.5230195407313964)]
[np.float64(0.5153133786232648), np.float64(0.8153646920821799),
np.float64(0.21425423974065558), np.float64(0.38189768467369134)]
[np.float64(0.5437598552020154), np.float64(0.9112448342915718),
np.float64(0.2325001490145239), np.float64(0.5216526317547565)]
[np.float64(0.550056946262832), np.float64(0.6881078683113983),
np.float64(0.20131810361017977), np.float64(0.4209082106732165)]
```

```
[np.float64(0.42208106874816853), np.float64(0.6881078683113983),
np.float64(0.18648179657124278), np.float64(0.4209082106732165)]
[np.float64(0.5252978276017184), np.float64(0.5477078274217113),
np.float64(0.16186101321783258), np.float64(0.5269119316680196)]
[np.float64(0.22406173998598164), np.float64(0.9605192178192339),
np.float64(0.06319073974654746), np.float64(0.15241623272983265)]
[np.float64(0.22406173998598164), np.float64(0.41117733955798913),
np.float64(0.2325001490145239), np.float64(0.4852191100646917)]
[np.float64(0.06949959924874023), np.float64(0.8953883092545903),
np.float64(0.05722734947336401), np.float64(0.15241623272983265)]
[np.float64(0.31288342509729494), np.float64(0.41117733955798913),
np.float64(0.18995792889043503), np.float64(0.48904343555449875)]
[np.float64(0.28535038299063187), np.float64(0.7918552521983719),
np.float64(0.15330075552164707), np.float64(0.34051002673066694)]
[np.float64(0.359243321521318), np.float64(0.5196105597681255),
np.float64(0.06952696936704778), np.float64(0.15241623272983265)]
[np.float64(0.550056946262832), np.float64(0.8153646920821799),
np.float64(0.20131810361017977), np.float64(0.21737923326444122)]
[np.float64(0.3776222618737363), np.float64(0.9605192178192339),
np.float64(0.1941040565258442), np.float64(0.38189768467369134)]
[np.float64(0.4846207658370162), np.float64(0.5196105597681255),
np.float64(0.05722734947336401), np.float64(0.15241623272983265)]
[np.float64(0.49781718112370765), np.float64(0.6604926430655316),
np.float64(0.17852059650837715), np.float64(0.2569007773342541)]
[np.float64(0.05847559072932759), np.float64(0.8153646920821799),
np.float64(0.11902905254890522), np.float64(0.4209082106732165)]
[np.float64(0.550056946262832), np.float64(0.6881078683113983),
np.float64(0.23613779832759446), np.float64(0.510177203090685)]
[np.float64(0.4628070552230153), np.float64(0.7918552521983719),
np.float64(0.12748254603026224), np.float64(0.1807565062254383)]
[np.float64(0.31288342509729494), np.float64(0.17808802830464457),
np.float64(0.07622601409392568), np.float64(0.15241623272983265)]
[np.float64(0.5375348536877407), np.float64(0.17808802830464457),
np.float64(0.08327280126275306), np.float64(0.18603539180084086)]
[np.float64(0.5375348536877407), np.float64(0.0043094204234492438),
np.float64(0.06952696936704778), np.float64(0.2231429525232999)]
[np.float64(0.4846207658370162), np.float64(0.5477078274217113),
np.float64(0.23519160623165483), np.float64(0.510177203090685)]
[np.float64(0.447093580707753), np.float64(0.6604926430655316),
np.float64(0.08327280126275306), np.float64(0.15241623272983265)]
```

As a test, take one datapoint from one of the other classes and see if the predicted probability changes.

Think a bit why the probability changes, what could affect the prediction?

## 1.5 Prepare Naive Bayes for binary classification

### 1.5.1 6) Prepare the data for inference

Before we train and test the naive bayes for multiple classes, we should get our data in order.

Similar to how we did previously, we should now split two classes into a train and test set, you may choose which two classes freely.

```python
# TODO: Split two classes into train and test sets.

# already done in my class flower

# TODO: Sepearte the features and lables for both the train and test set.

# already done in my class flower


# Class A : iris_setosa

# Class B : iris_versicolor
```

### 1.5.2 7) Class A vs Class B for binary classification

**Note:** You might need to go back and forth a bit in the following cells during your implementation of your code.

We have to get the probability from two sets of classes and compare the two probabilities in order to make a propper prediction.

Here we will implement two functions to make this possible. We seperate these functions to make the implementation of the ROC-curve easier later on.

**Function 1: naive_bayes_prediction** * A function that returns the probabilities for each class the model for a single datapoint.

**Function 2: probabilities_to_prediction** * A function that takes in probabilities and returns a prediction.

```python
def naive_bayes_prediction(feature_stats, data_point):
    probabilities = {}

    for class_label, stats in feature_stats.items():
        mean, stdev, prior = stats['mean'], stats['stdev'], stats['prior']

        # Probability for this class
        # print("prior", prior)
        prob = float(prior)
        for x, m, s in zip(data_point, mean, stdev):
            # print(f"x: {x}, mean: {m}, std: {s}, prob: {prob}")
            prob *= gaussian_probability_density_function(x, m, s)
```

```
        probabilities[class_label] = prob

    return probabilities
```

```
[133]: def probabilities_to_prediction(probabilities):
           class_prediction = max(probabilities, key=probabilities.get)
           return class_prediction
```

To test the function we need the feature metrics from the classes we choose.

**Note:** Choose the correct train/test set and the correct feature/label split!

```
[134]: # TODO: Get the feature metrics for the classes.
```

Now we should have implemented all the neccessary parts to train a naive bayes algorithm and do inference on it. Implement a small test workflow for two of your chosen classes.

```
[135]: def evaluate_and_get_probs(item, feature_stats, actual_class_name,␣
       ↪positive_class_name):
           correct = 0
           total = len(item.test.features)

           prediction_probabilities = []
           test_labels = []

           for x, true_label in zip(item.test.features, item.test.labels):
               probs = naive_bayes_prediction(feature_stats, x)
               prediction = probabilities_to_prediction(probs)

               is_correct = (prediction == actual_class_name)
               correct += int(is_correct)

               print(f"Features: {x}, Prediction = {prediction}, Actual =␣
       ↪{actual_class_name}")

               p_pos = probs[positive_class_name]
               prediction_probabilities.append(p_pos)

               if actual_class_name == positive_class_name:
                   test_labels.append(1)
               else:
                   test_labels.append(0)

           accuracy_test = correct / total
           print(f"\nAccuracy on {actual_class_name} test samples = {accuracy_test:.
       ↪2f}\n")
```

```
        return prediction_probabilities, test_labels
```

```
[136]: feature_stats = {
           "setosa": {
               'mean': iris_setosa.train.means,
               'stdev': iris_setosa.train.std_devs,
               'prior': iris_setosa.train.prior
           },
           "versicolor": {
               'mean': iris_versicolor.train.means,
               'stdev': iris_versicolor.train.std_devs,
               'prior': iris_versicolor.train.prior
           }
       }

       pred_probs_setosa, labels_setosa = evaluate_and_get_probs(
           iris_setosa,
           feature_stats,
           actual_class_name="setosa",
           positive_class_name="setosa"
       )

       pred_probs_versi, labels_versi = evaluate_and_get_probs(
           iris_versicolor,
           feature_stats,
           actual_class_name="versicolor",
           positive_class_name="versicolor"
       )

       all_pred_probs = pred_probs_setosa + pred_probs_versi
       all_labels = labels_setosa + labels_versi
```

```
Features: [5.1 3.5 1.4 0.2], Prediction = setosa, Actual = setosa
Features: [5.1 3.8 1.5 0.3], Prediction = setosa, Actual = setosa
Features: [6.8 3.2 5.9 2.3], Prediction = versicolor, Actual = setosa
Features: [5.8 2.6 4.  1.2], Prediction = versicolor, Actual = setosa
Features: [5.  3.2 1.2 0.2], Prediction = setosa, Actual = setosa
Features: [6.9 3.1 4.9 1.5], Prediction = versicolor, Actual = setosa
Features: [5.  3.3 1.4 0.2], Prediction = setosa, Actual = setosa
Features: [4.4 3.  1.3 0.2], Prediction = setosa, Actual = setosa
Features: [4.9 3.1 1.5 0.1], Prediction = setosa, Actual = setosa
Features: [6.3 2.5 5.  1.9], Prediction = versicolor, Actual = setosa

Accuracy on setosa test samples = 0.60

Features: [7.1 3.  5.9 2.1], Prediction = versicolor, Actual = versicolor
Features: [6.2 2.9 4.3 1.3], Prediction = versicolor, Actual = versicolor
Features: [6.7 3.  5.2 2.3], Prediction = versicolor, Actual = versicolor
```

```
Features: [5.9 3.  5.1 1.8], Prediction = versicolor, Actual = versicolor
Features: [4.9 3.1 1.5 0.1], Prediction = setosa, Actual = versicolor
Features: [7.  3.2 4.7 1.4], Prediction = setosa, Actual = versicolor
Features: [6.1 3.  4.9 1.8], Prediction = versicolor, Actual = versicolor
Features: [5.  2.3 3.3 1. ], Prediction = versicolor, Actual = versicolor
Features: [5.6 2.7 4.2 1.3], Prediction = versicolor, Actual = versicolor
Features: [4.8 3.  1.4 0.3], Prediction = setosa, Actual = versicolor

Accuracy on versicolor test samples = 0.70
```

## 1.6  ROC-curve

A ROC curve, or *Receiver Operating Characteristic curve*, is a graphical plot that illustrates the performance of a binary classifier such as our Naive Bayes model.

More info can be found in the course material and here: https://en.wikipedia.org/wiki/Receiver_operating_characteristic

Another good illustration by Google can be found here: https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc

Now that we have a prediction model, we would want to try it out and test it using a ROC-curve.

### 1.6.1  8) True Positive Rate (TPR) and False Positive Rate (FPR)

From our prediction function we get probabilities, and for prediction purposes we have just predicted the one with the highest probability.

To plot a ROC-curve, we need the TPR and FPR for the binary classification. We will implement this here.

**Note 1:** The threshold is is a value that goes from 0 to 1.

**Note 2:** One of the two classes will be seen as "the positive class" (prediction over the threshold) and the other as "the negative class" (prediction under the threshold).

**Note 3:** The threshold stepsize will decide the size of the returned TPR/FPR list. A value of 0.1 will give 10 elements (0 to 1 in increments of 0.1)

```
[137]: # Stepsize demonstration
       print("Python list:", [x/10 for x in range(0,10,1)])

       # Stepsize demonstration with numpy:
       print("Numpy linspace:", np.linspace(0,1,11))
       print("Numpy linspace (no endpoint):", np.linspace(0,1,10,endpoint=False))
```

```
Python list: [0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
Numpy linspace: [0.  0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1. ]
Numpy linspace (no endpoint): [0.  0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9]
```

```
[138]: def TPR_and_FPR(prediction_probabilities, test_labels, threshold_stepsize=0.1):
            thresholds = np.linspace(0, 1, int(1/threshold_stepsize) + 1)
            TPR_list = []
            FPR_list = []

            prediction_probabilities = np.array(prediction_probabilities, dtype=float)
            test_labels = np.array(test_labels, dtype=int)

            for threshold in thresholds:
                predicted_labels = (prediction_probabilities >= threshold).astype(int)

                TP = np.sum((predicted_labels == 1) & (test_labels == 1))
                FP = np.sum((predicted_labels == 1) & (test_labels == 0))
                TN = np.sum((predicted_labels == 0) & (test_labels == 0))
                FN = np.sum((predicted_labels == 0) & (test_labels == 1))

                TPR = TP / (TP + FN) if (TP + FN) > 0 else 0.0
                FPR = FP / (FP + TN) if (FP + TN) > 0 else 0.0

                TPR_list.append(TPR)
                FPR_list.append(FPR)

            return TPR_list, FPR_list
```

```
[139]: # TODO: Test the "TPR_and_FPR" function on the model you have created␣
       ↪previously.
       TPR, FPR = TPR_and_FPR(all_pred_probs, all_labels, threshold_stepsize=0.1)
       print("TPR =", TPR)
       print("FPR =", FPR)
```

```
TPR = [np.float64(1.0), np.float64(0.0), np.float64(0.0), np.float64(0.0),
np.float64(0.0), np.float64(0.0), np.float64(0.0), np.float64(0.0),
np.float64(0.0), np.float64(0.0), np.float64(0.0)]
FPR = [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]
```

How does the values change if you change the threshold stepsize?

How does the values change if you change the classes you compare?

### 1.6.2 9) Plot the TPR and FPR

To better see what is going on, we can plot the TPR and FPR. We can also calculate the Area Under the ROC Curve (AUC or AUROC) at the same time.

```
[140]: def plot_ROC(TPR, FPR):

           auc_score = np.trapz(TPR, x=FPR)

           plt.figure(figsize=(6, 6))
```

```python
        plt.plot(FPR_sorted, TPR_sorted, label=f"ROC (AUC={auc_score:.3f})")
        plt.plot([0, 1], [0, 1], 'r--')
        plt.xlabel("FPR")
        plt.ylabel("TPR")
        plt.title("ROC Curve")
        plt.legend(loc="lower right")
        plt.grid(True)
        plt.show()

        print(f"AUC : {auc_score:.3f}")
```

```python
[141]:  # TODO: Test the plotting function on the TPR and FPR you just calculated.
        TPR, FPR = TPR_and_FPR(all_pred_probs, all_labels, threshold_stepsize=0.1)

        pairs = sorted(zip(FPR, TPR), key=lambda x: x[0])
        FPR_sorted, TPR_sorted = zip(*pairs)


        plot_ROC(TPR_sorted, FPR_sorted)
```
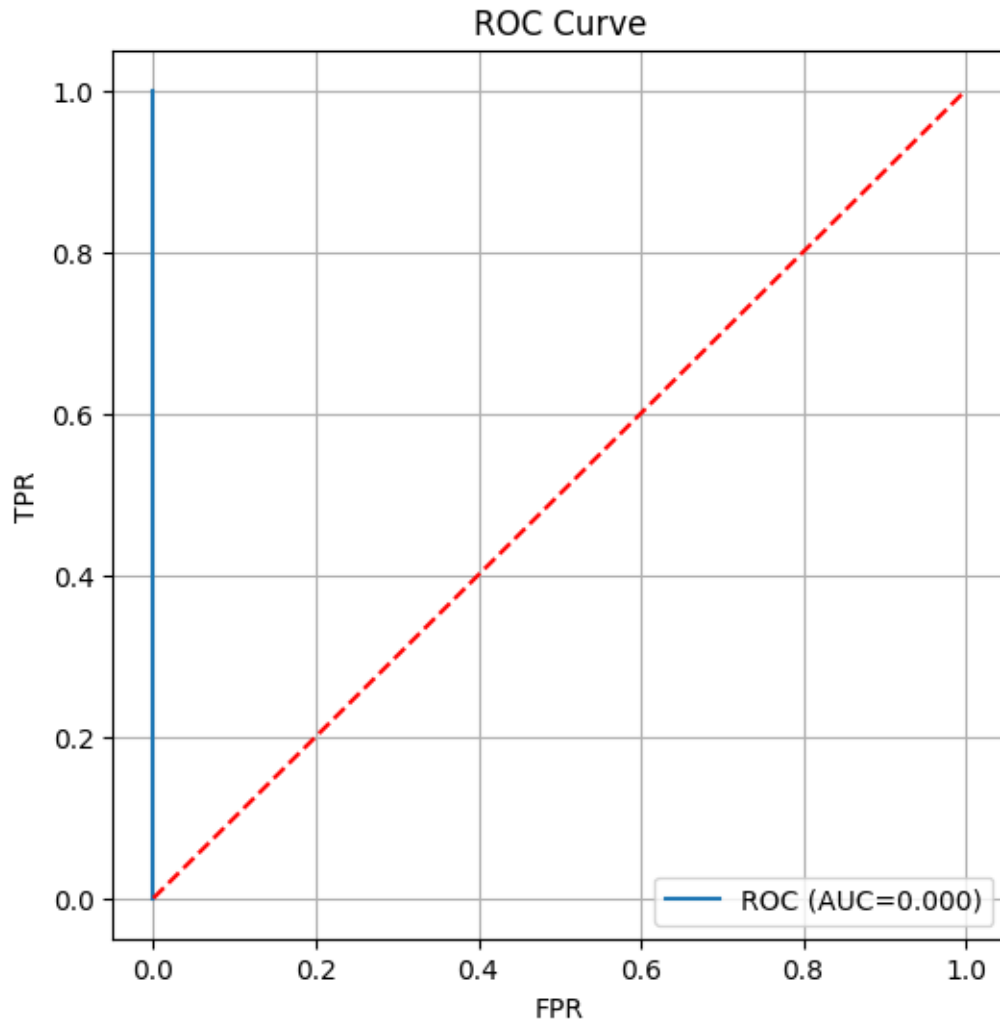
/var/folders/yx/0678sjj54074f8593p3dv0cc0000gn/T/ipykernel_7124/1290740080.py:3:
DeprecationWarning: `trapz` is deprecated. Use `trapezoid` instead, or one of
the numerical integration functions in `scipy.integrate`.
  auc_score = np.trapz(TPR, x=FPR)

```
AUC : 0.000
```

## 1.7 Cross-validation

The final task is to take everything you have implemented so far and apply it in a cross-validation loop.

**Note 1:** To better reflect a real scenarios, you should shuffle the data before doing cross-validation.

**Note 2:** When using cross-validation, the interesting thing is the mean performance (mean AUC, mean accuracy, mean ROC-curve).

**Note 3:** This part is a bit more free in terms of implementation, but make sure to use some of the previously implemented functions.

### 1.7.1  10) Cross-validation loop

```python
[142]: def cross_validation(features, labels, folds=5, threshold_stepsize=0.1):
           features = np.array(features, dtype=float)
           labels = np.array(labels, dtype=int)
           N = len(labels)

           indices = np.arange(N)
           np.random.shuffle(indices)

           fold_size = N // folds

           all_TPR = []
           all_FPR = []

           for i in range(folds):
               start = i * fold_size
               end = (i+1)*fold_size if (i < folds-1) else N

               test_indices = indices[start:end]

               train_indices = np.concatenate((indices[:start], indices[end:]))

               X_train = features[train_indices]
               y_train = labels[train_indices]
               X_test  = features[test_indices]
               y_test  = labels[test_indices]

               X_train_pos = X_train[y_train == 1]
               X_train_neg = X_train[y_train == 0]

               means_pos, stdevs_pos, prior_pos =␣
        ↪calculate_feature_statistics(X_train_pos, X_train)
               means_neg, stdevs_neg, prior_neg =␣
        ↪calculate_feature_statistics(X_train_neg, X_train)

               feature_stats_fold = {
                   1: {
                       'mean':  means_pos,
                       'stdev': stdevs_pos,
                       'prior': prior_pos
                   },
                   0: {
                       'mean':  means_neg,
                       'stdev': stdevs_neg,
                       'prior': prior_neg
                   }
```

```
        }

        prediction_probabilities = []
        for x_vec in X_test:
            probs = naive_bayes_prediction(feature_stats_fold, x_vec)
            prob_pos = probs[1]
            prediction_probabilities.append(prob_pos)

        fold_TPR, fold_FPR = TPR_and_FPR(prediction_probabilities, y_test,␣
  ↪threshold_stepsize)
        all_TPR.append(fold_TPR)
        all_FPR.append(fold_FPR)

    all_TPR = np.array(all_TPR)
    all_FPR = np.array(all_FPR)

    mean_TPR = np.mean(all_TPR, axis=0)
    mean_FPR = np.mean(all_FPR, axis=0)

    plot_ROC(mean_TPR, mean_FPR)
```

### 1.7.2 11) 10-fold Cross-validation on all classes

Test the "cross_validation" function on all the classes against eachother using 10 folds.

- Iris-setosa vs Iris-versicolor
- Iris-setosa vs Iris-virginica
- Iris-versicolor vs Iris-virginica

```
[143]: # TODO: Implement and test cross-validation function on all classes.


       features_setosa = iris_setosa.train.features
       features_versi = iris_versicolor.train.features

       labels_setosa = np.ones(len(features_setosa), dtype=int)
       labels_versi = np.zeros(len(features_versi), dtype=int)

       all_features = np.vstack([features_setosa, features_versi])
       all_labels = np.concatenate([labels_setosa, labels_versi])

       print("Cross-validation Setosa vs Versicolor")
       cross_validation(all_features, all_labels, folds=10, threshold_stepsize=0.1)
```
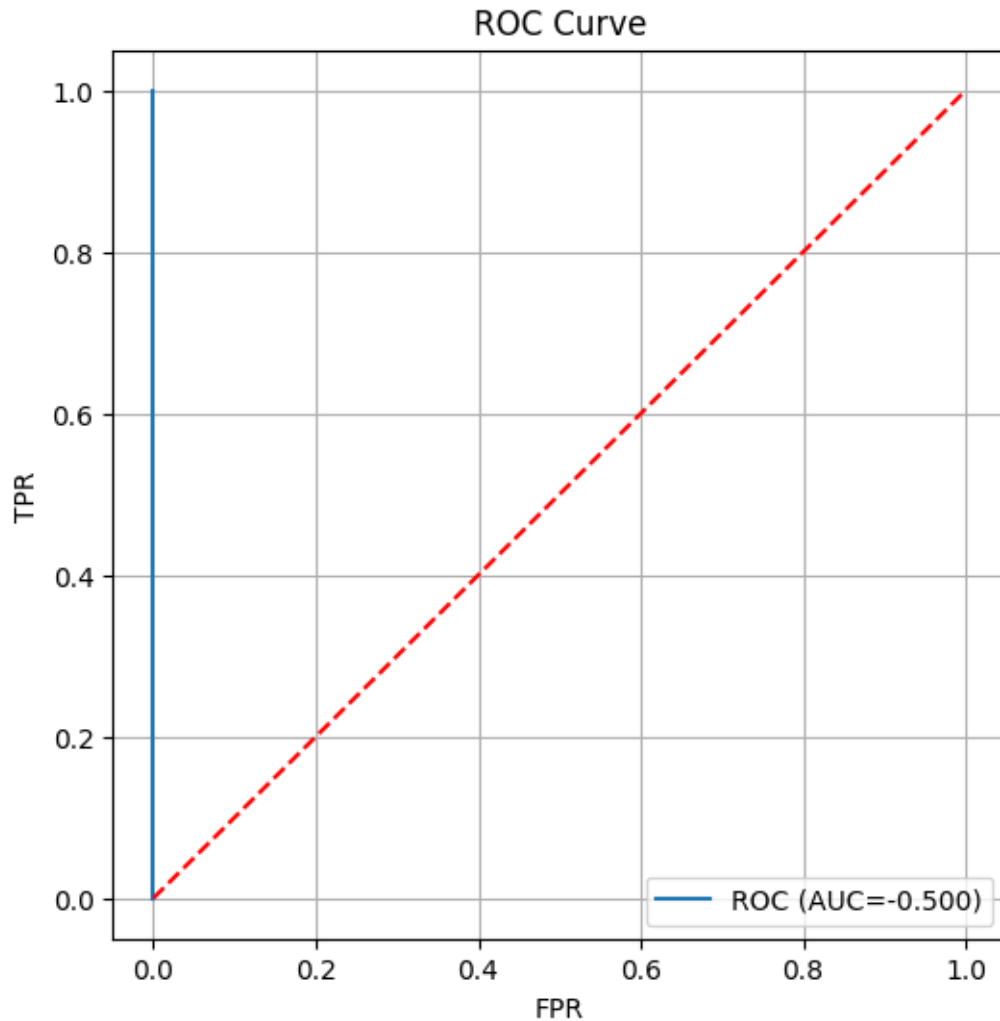
Cross-validation Setosa vs Versicolor

/var/folders/yx/0678sjj54074f8593p3dv0cc0000gn/T/ipykernel_7124/1290740080.py:3:
DeprecationWarning: `trapz` is deprecated. Use `trapezoid` instead, or one of
the numerical integration functions in `scipy.integrate`.

```
auc_score = np.trapz(TPR, x=FPR)
```

## ROC Curve



AUC : -0.500

```
[144]: # TODO: Implement and test cross-validation function on all classes.
       features_setosa = iris_setosa.train.features
       features_virgi = iris_virginica.train.features

       labels_setosa = np.ones(len(features_setosa), dtype=int)
       labels_virgi = np.zeros(len(features_virgi), dtype=int)

       all_features = np.vstack([features_setosa, features_virgi])
       all_labels = np.concatenate([labels_setosa, labels_virgi])

       print("Cross-validation Setosa vs Virginica")
```
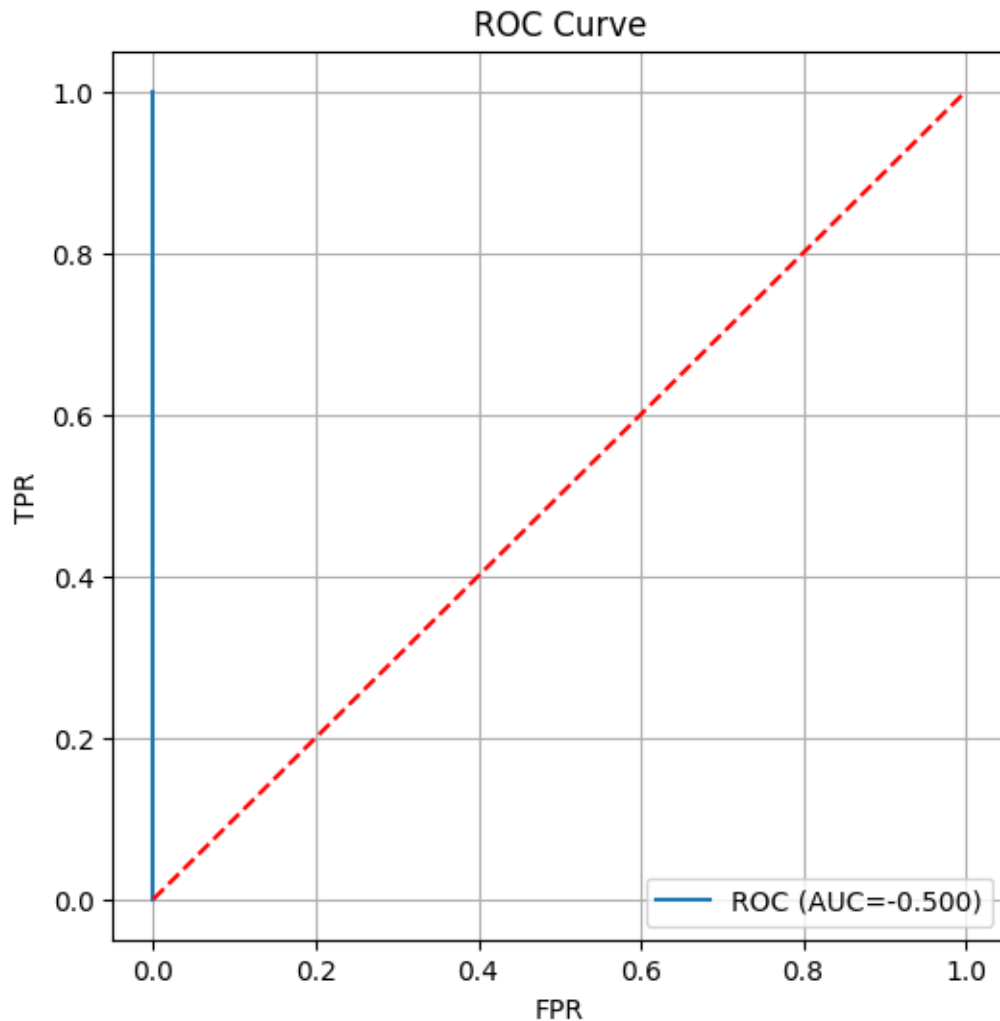
```
cross_validation(all_features, all_labels, folds=10, threshold_stepsize=0.1)
```

Cross-validation Setosa vs Virginica

/var/folders/yx/0678sjj54074f8593p3dv0cc0000gn/T/ipykernel_7124/1290740080.py:3:
DeprecationWarning: `trapz` is deprecated. Use `trapezoid` instead, or one of
the numerical integration functions in `scipy.integrate`.
  auc_score = np.trapz(TPR, x=FPR)



```
AUC : -0.500
```

[145]:
```python
# TODO: Implement and test cross-validation function on all classes.
features_versi = iris_versicolor.train.features
features_virgi = iris_virginica.train.features

labels_versi = np.ones(len(features_versi), dtype=int)
```

```
labels_virgi = np.zeros(len(features_virgi), dtype=int)

all_features = np.vstack([features_versi, features_virgi])
all_labels = np.concatenate([labels_versi, labels_virgi])

print("Cross-validation Versicolor vs Virginica")
cross_validation(all_features, all_labels, folds=10, threshold_stepsize=0.1)
```
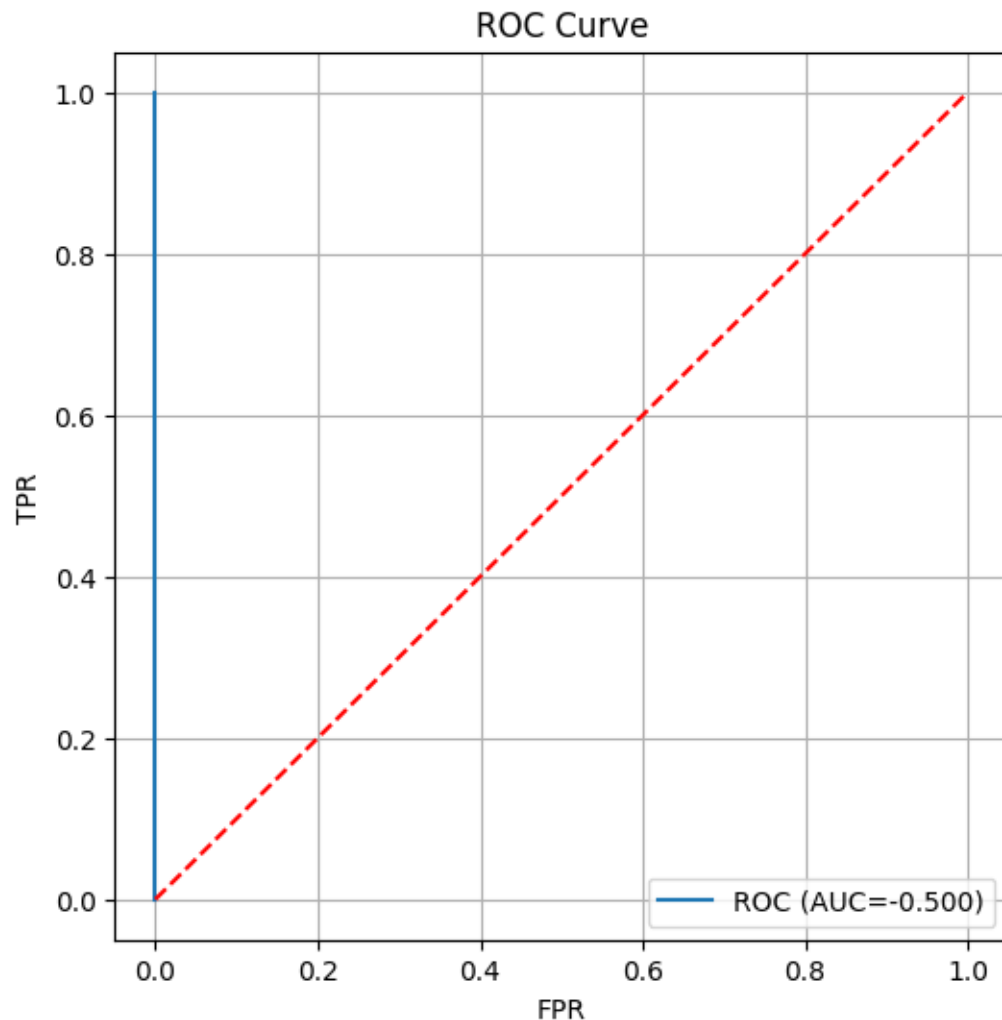
Cross-validation Versicolor vs Virginica

/var/folders/yx/0678sjj54074f8593p3dv0cc0000gn/T/ipykernel_7124/1290740080.py:3:
DeprecationWarning: `trapz` is deprecated. Use `trapezoid` instead, or one of
the numerical integration functions in `scipy.integrate`.
  auc_score = np.trapz(TPR, x=FPR)



AUC : -0.500

## 2 Questions for examination:

In addition to completing the assignment with all its tasks, you should also prepare to answer the following questions:

1) Why is it called "naive bayes"?

2) What are some downsides of the naive bayes learning algorithm?

3) When using ROC-curves, what is the theoretical best and worst result you can get?

4) When using ROC-curves, in this assignment for example, is a higher threshold-stepsize always better?

5) When using cross-validation and ROC-curves, why is it important to take the correct mean values? What could go wrong?

## 3 Finished!

Was part of the setup incorrect? Did you spot any inconsistencies in the assignment? Could something improve?

If so, please write them and send via email and send it to:

- marcus.gullstrand@ju.se

Thank you!

```
[ ]:
```

```
[ ]:
```

```
[ ]:
```