

ITP 449

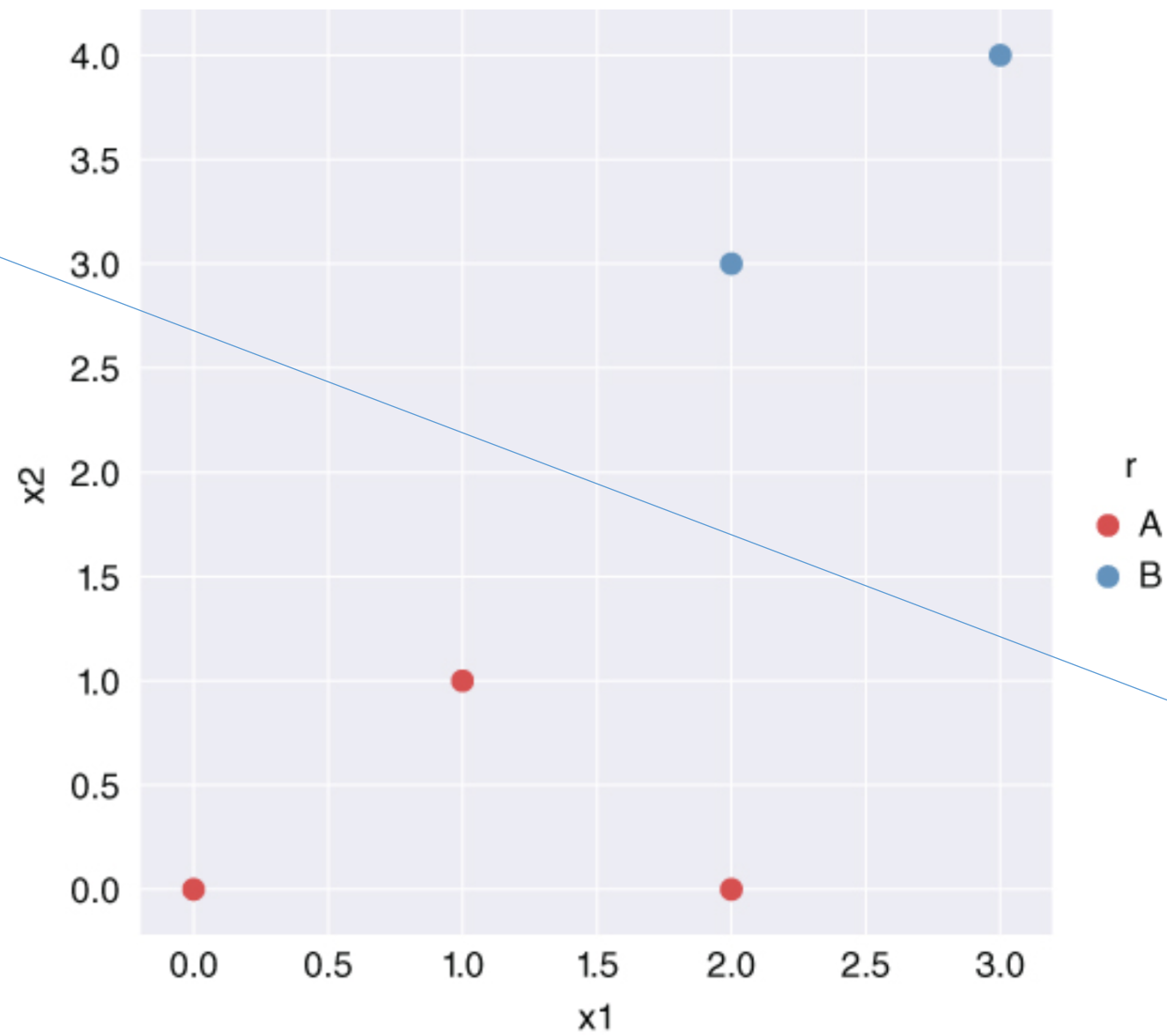
Support Vector Machines

Lecture 12



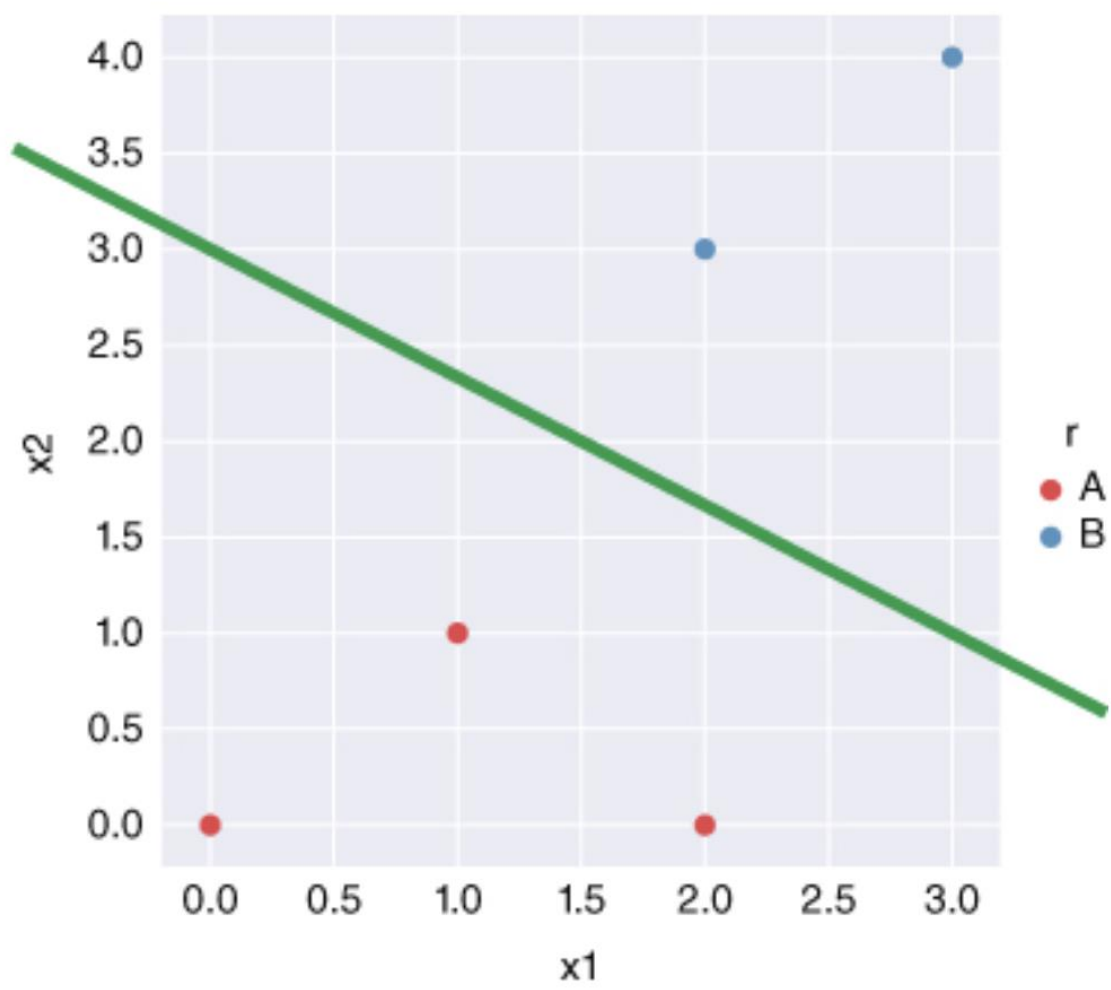
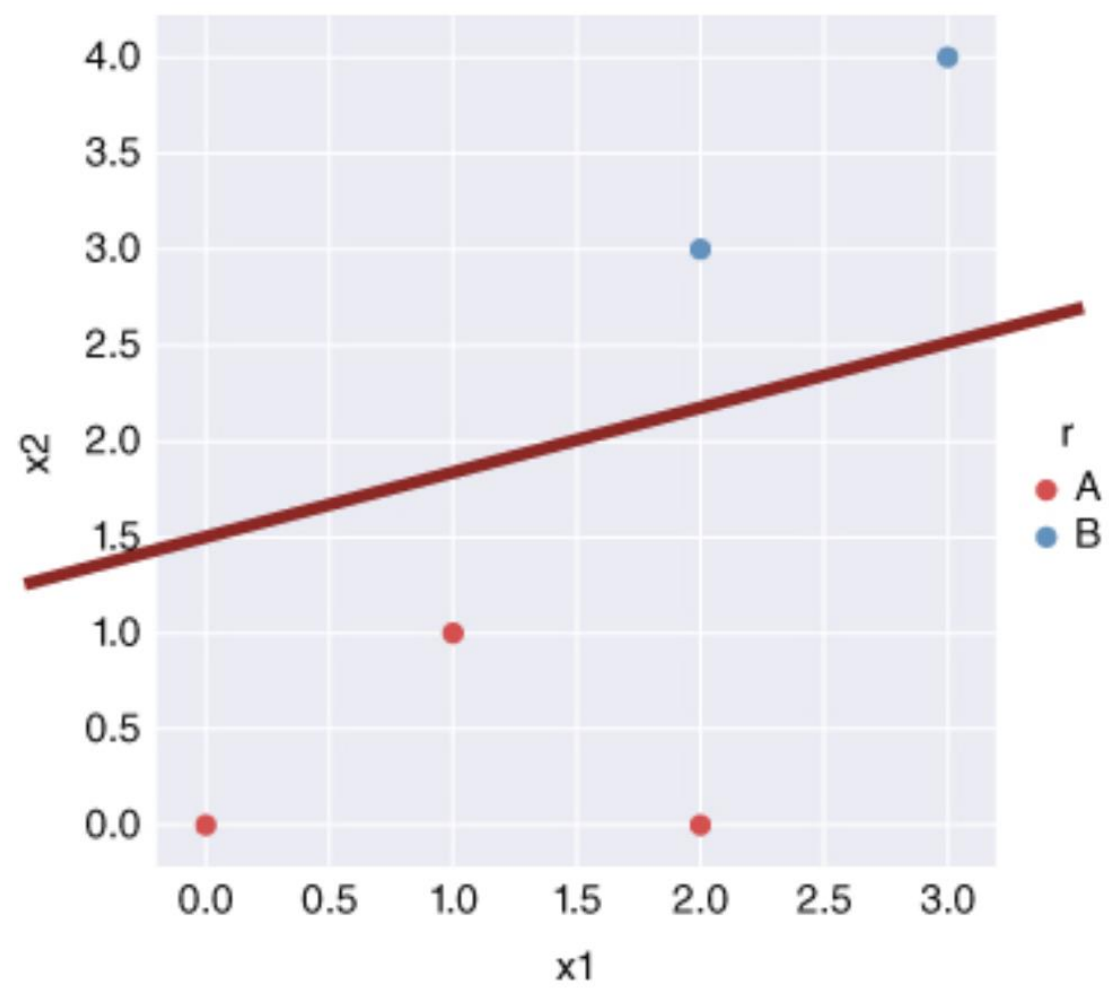
Support Vector Machines (SVMs)

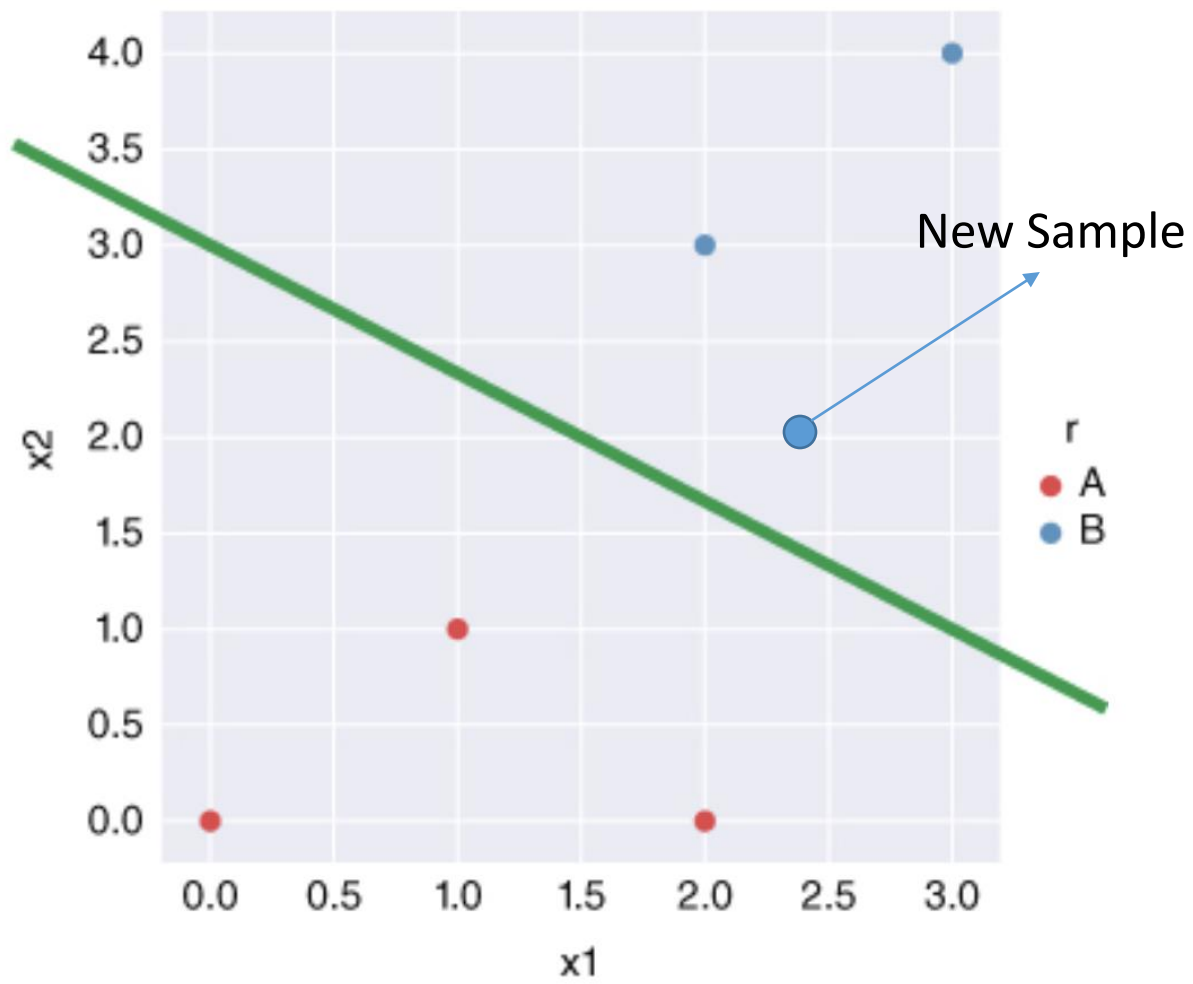
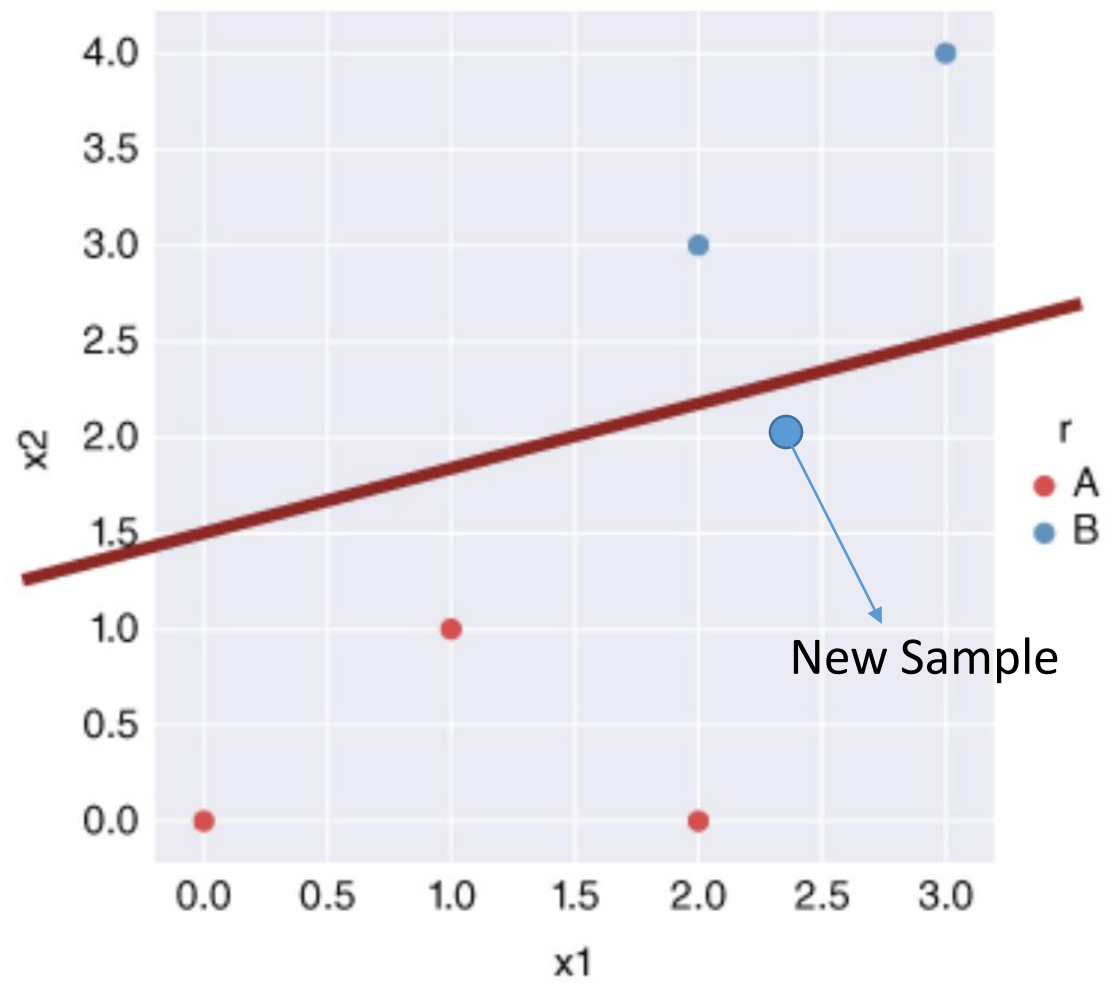
Used to predict categories and solve classification machine learning problems by defining a boundary between categories such that it maximizes the margin between classification groups.



Problem: Draw a line (hyperplane) on the chart so as to distinctly classify the data points into red and blue classes.

Many such lines are possible. Find one that yields the maximum margin. The highest distance between the points.



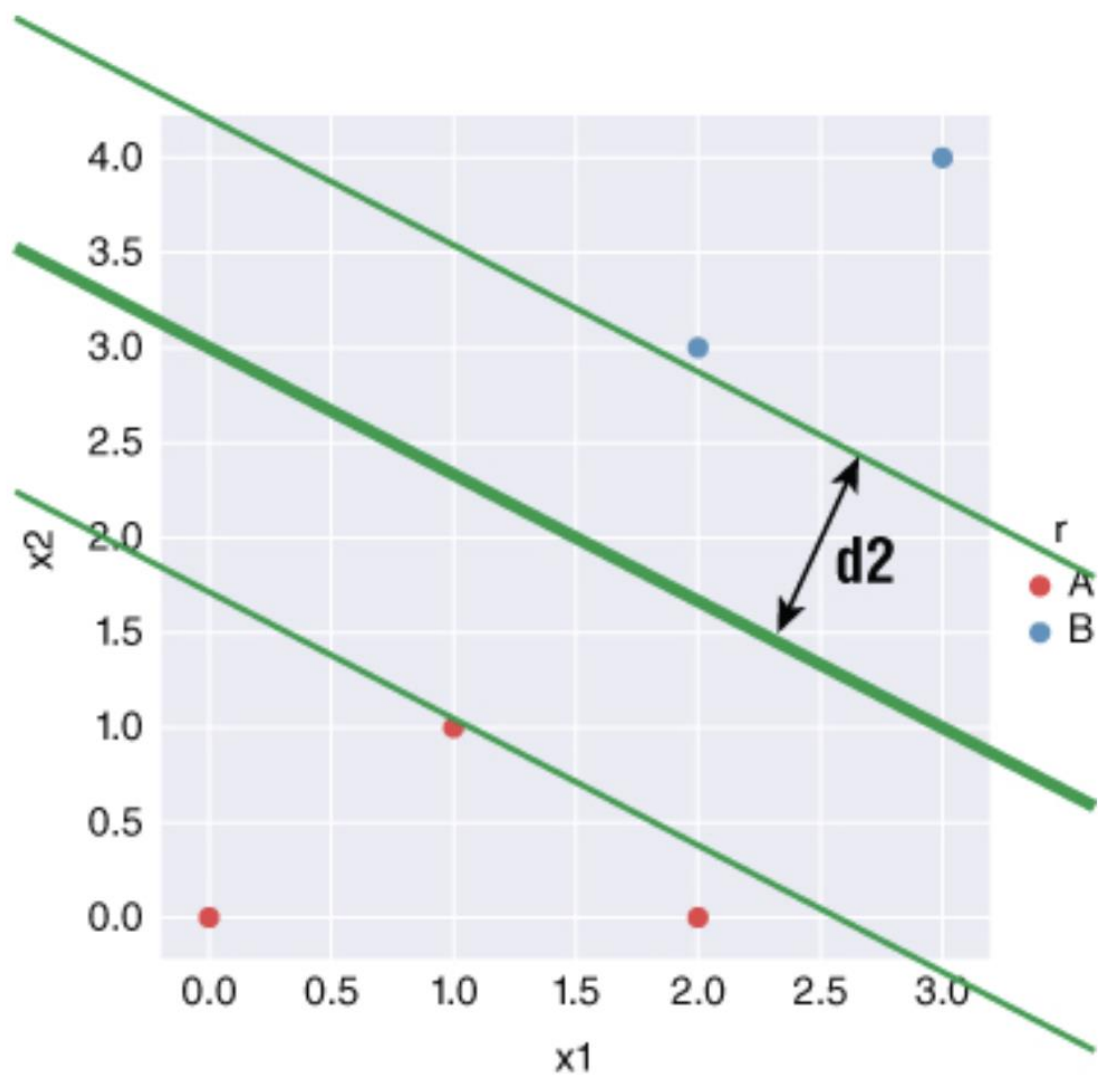
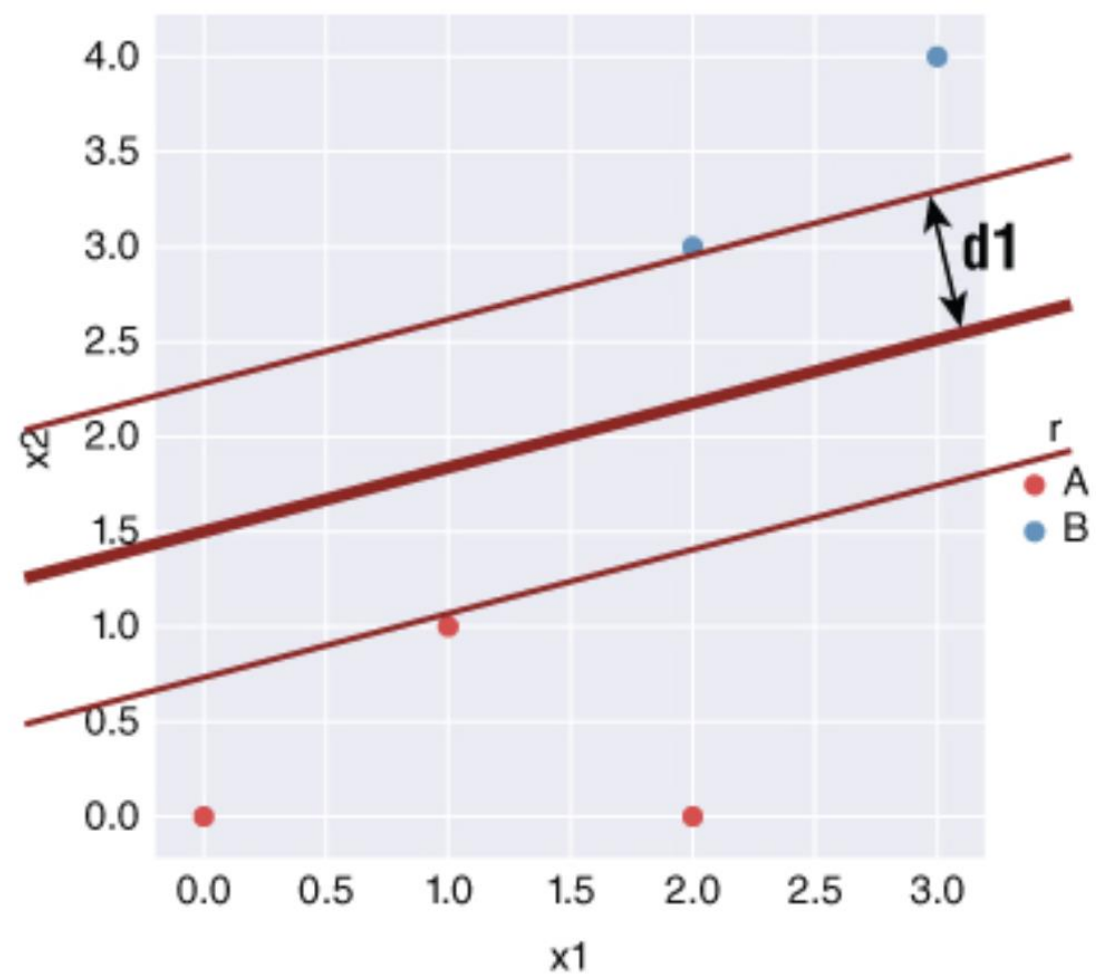


The bold line is the **hyperplane**. Also known as decision boundary.

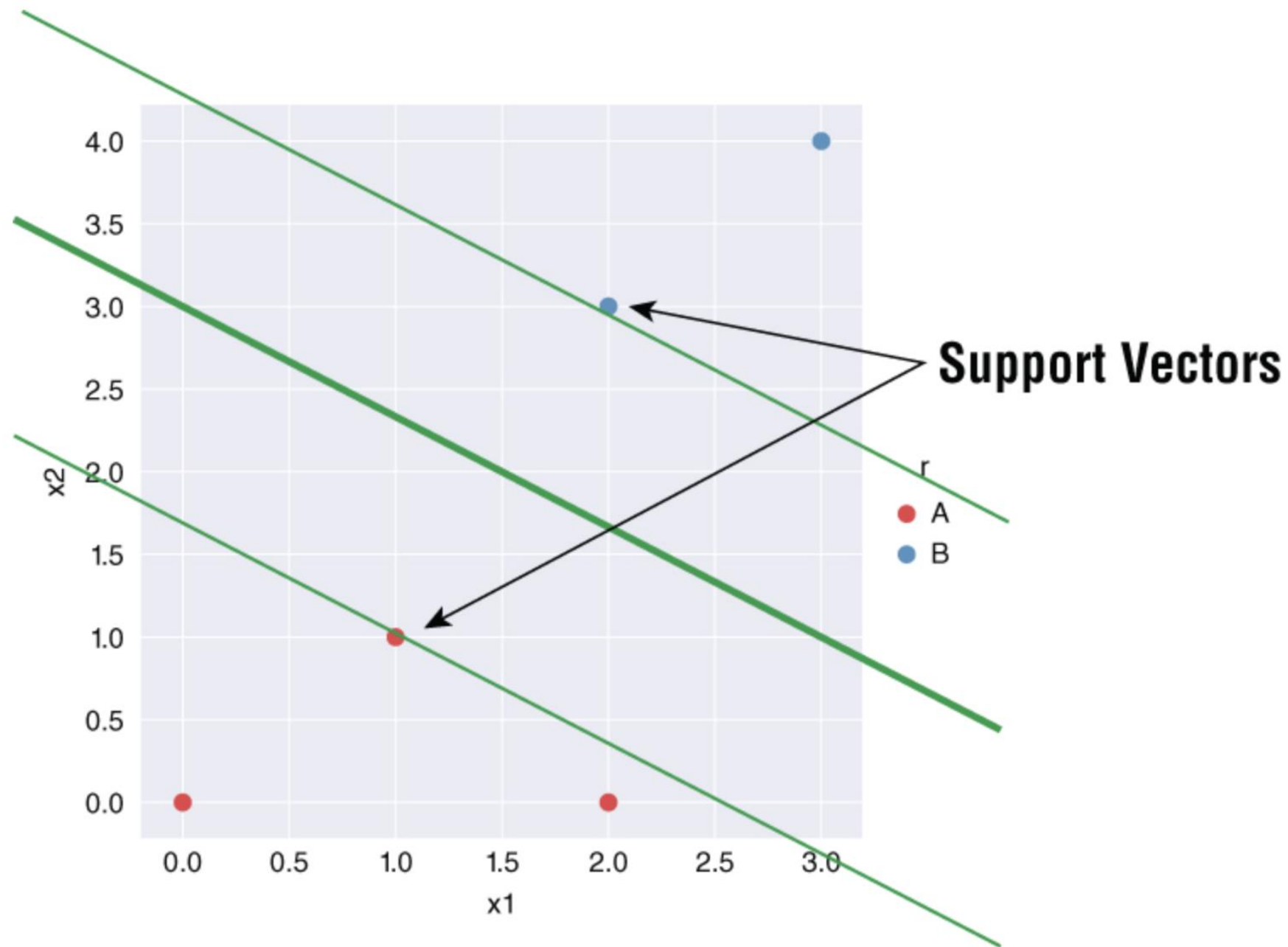
The reason to maximize the margin is that (new) points close to the decision boundary are not misclassified even if they are closer than the support vector. By having a large margin, you reduce this probability.

Support vectors are the points in each class that are close to the hyperplane.

Of course, there will be scenarios where it may not be possible to classify all points distinctly. There will be data points that end up in the opposite class. Instead of a straight line you could use a curved line to better separate the classes.



d2>d1



7: Structure

5 265 ^ v

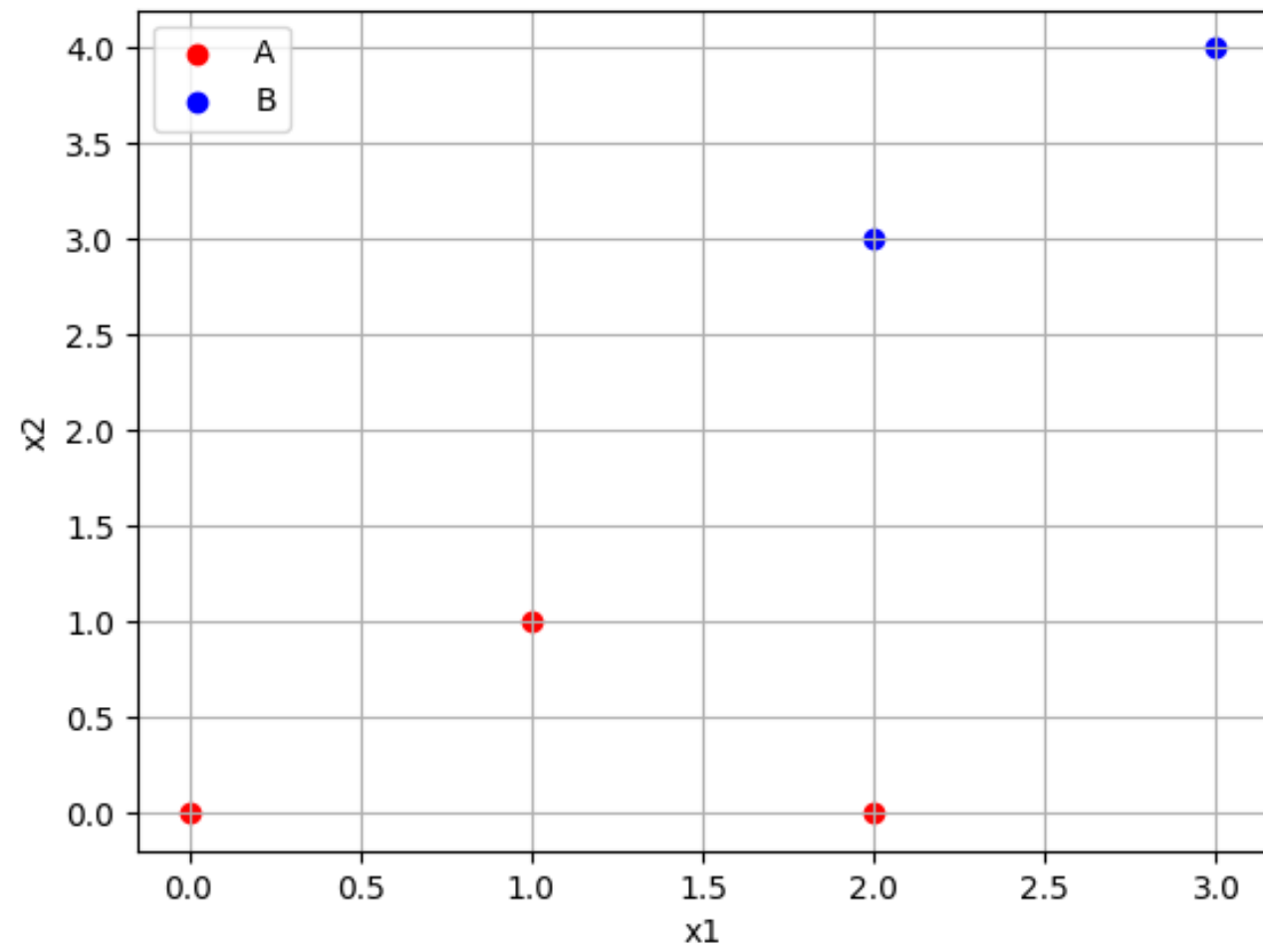
★ 2: Favorites

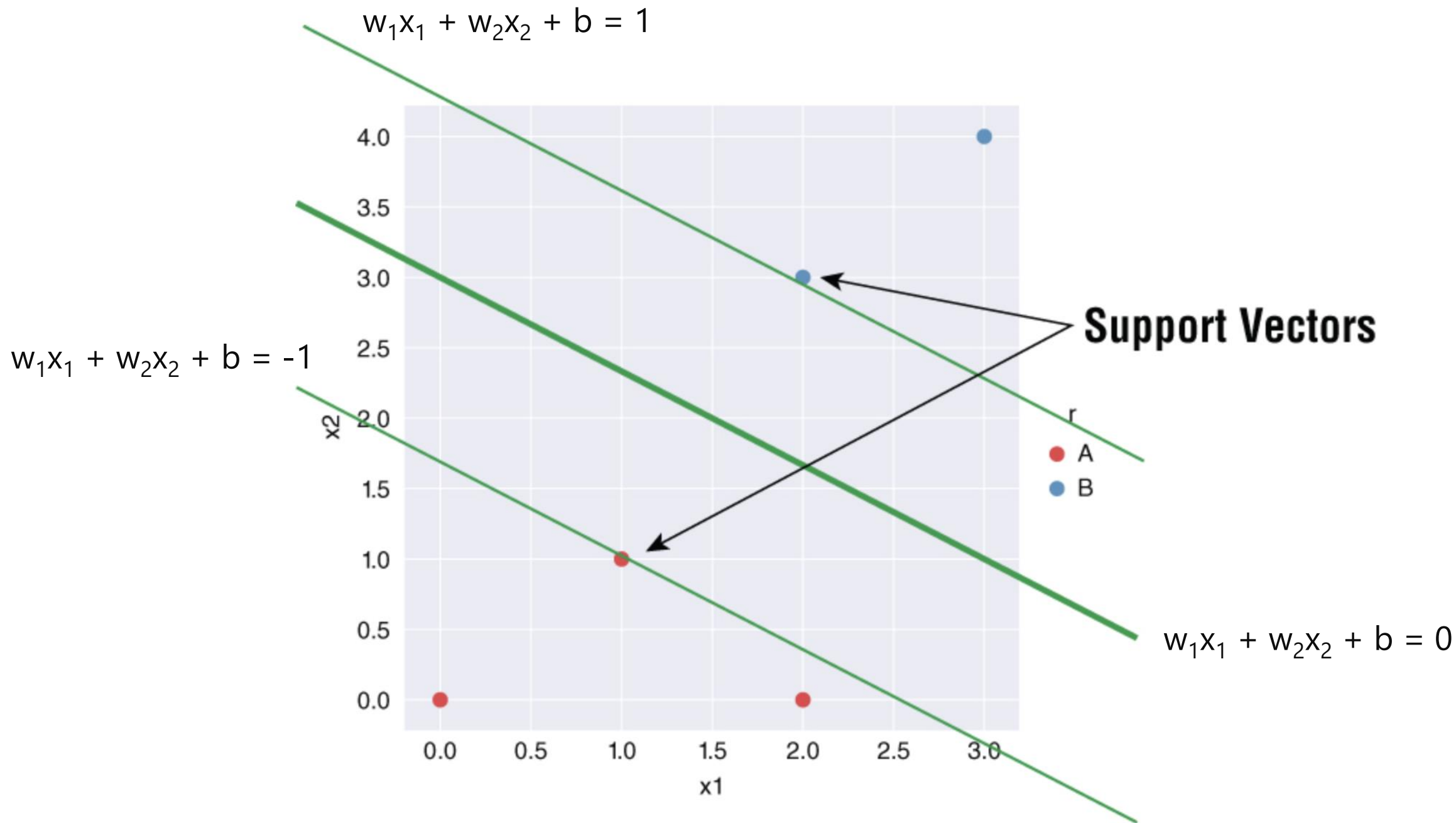
▶ 4: Run ≡ TODO ! 6: Problems ≡ Terminal 🐍 Python Console

Event Log



```
File Edit View Navigate Code Refactor Run Tools VCS Window Help ITP449_Fall2020 - in_class_coding.py
ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py in_class_coding
in_class_coding.py x
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 data = [[0, 0, 'A'], [1, 1, 'A'], [2, 3, 'B'], [2, 0, 'A'], [3, 4, 'B']]
5 df = pd.DataFrame(data)
6 df.columns = ['x1', 'x2', 'r']
7
8 plt.scatter(df.loc[df.r == 'A', 'x1'], df.loc[df.r == 'A', 'x2'], color='r', label='A')
9 plt.scatter(df.loc[df.r == 'B', 'x1'], df.loc[df.r == 'B', 'x2'], color='b', label='B')
10
11 plt.xlabel('x1')
12 plt.ylabel('x2')
13 plt.legend()
14 plt.grid()
15 plt.show()
16
Run: in_class_coding x
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"
Process finished with exit code 0
4: Run TODO 6: Problems Terminal Python Console Event Log
17:1 CRLF UTF-8 4 spaces Python 3.8 (Class)
```





1: Project in_class_coding.py

```

16
17 from sklearn import svm
18
19 model = svm.SVC(kernel='linear')
20
21 X = df[['x1', 'x2']]
22 y = df['r']
23
24 model.fit(X, y)
25
26 print('Weights w1 and w2:', model.coef_)
27 print('Bias:', model.intercept_)
28

```

6 265

SVC = Support Vector Classification
 Kernel = a technique to use higher dimensional space in order to make the classes separable (when they are not).
 Options: linear, polynomial, sigmoid etc.

Run: in_class_coding

```

C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"
Weights w1 and w2: [[0.4 0.8]]
Bias: [-2.2]

Process finished with exit code 0

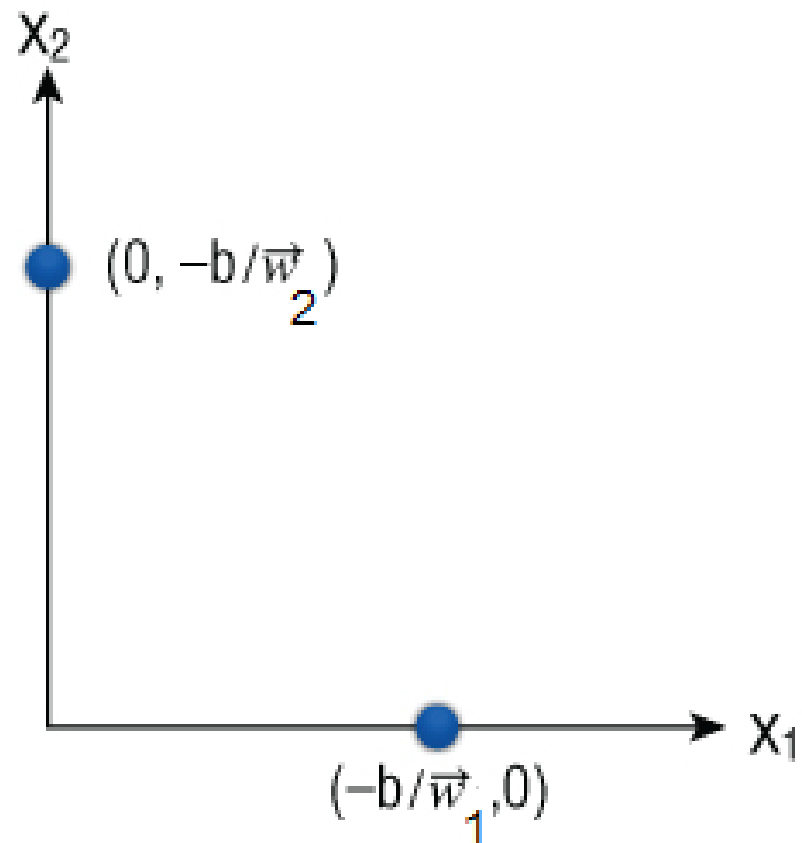
```

$$w_1x_1 + w_2x_2 + b = 0$$

$$w_1(0) + w_2x_2 + b = 0$$

$$w_2x_2 = -b$$

$$x_2 = -b / w_2$$



$$\text{slope} = (-b / w_2) / (b / w_1)$$

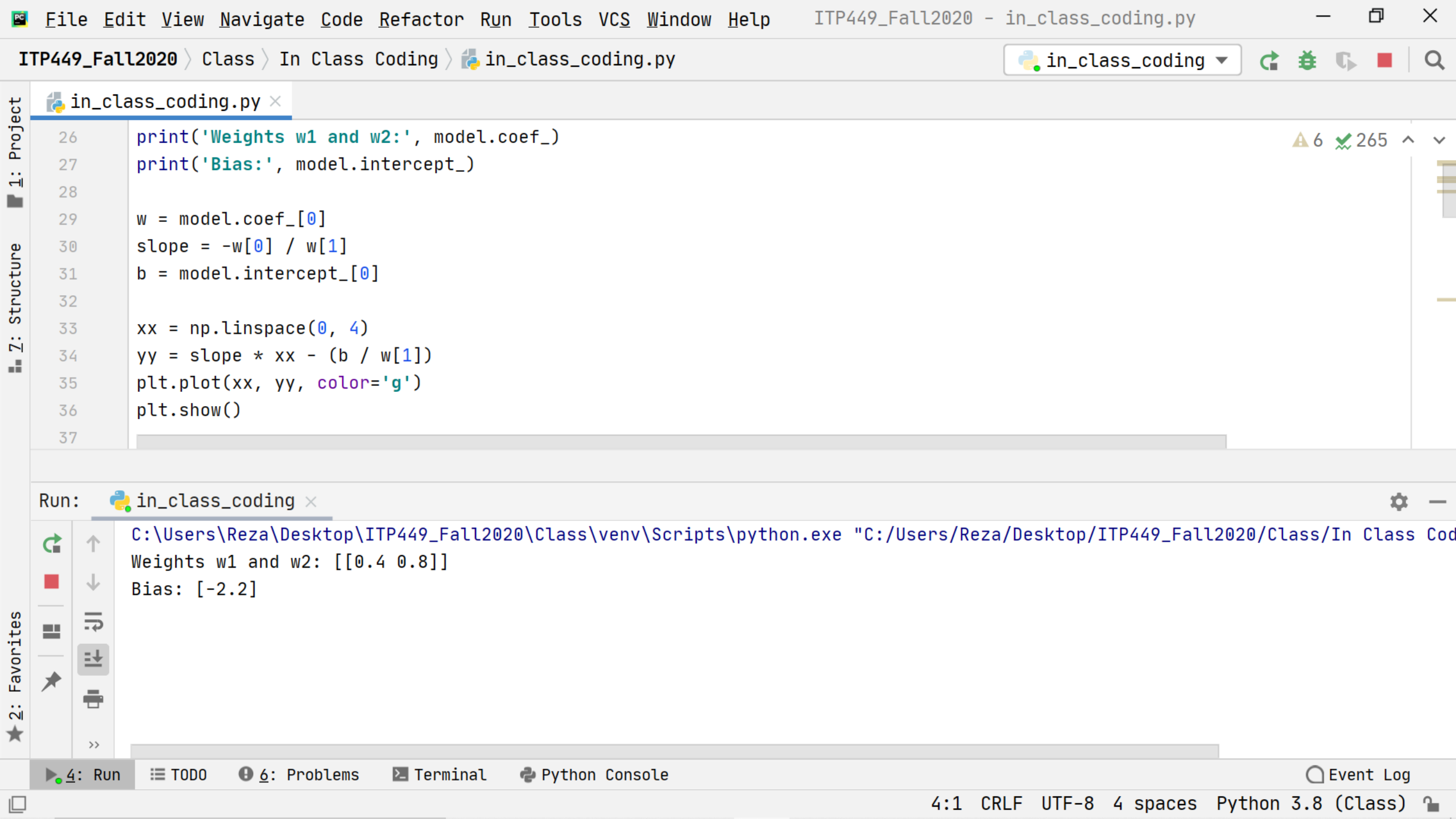
$$\text{slope} = - (w_1 / w_2)$$

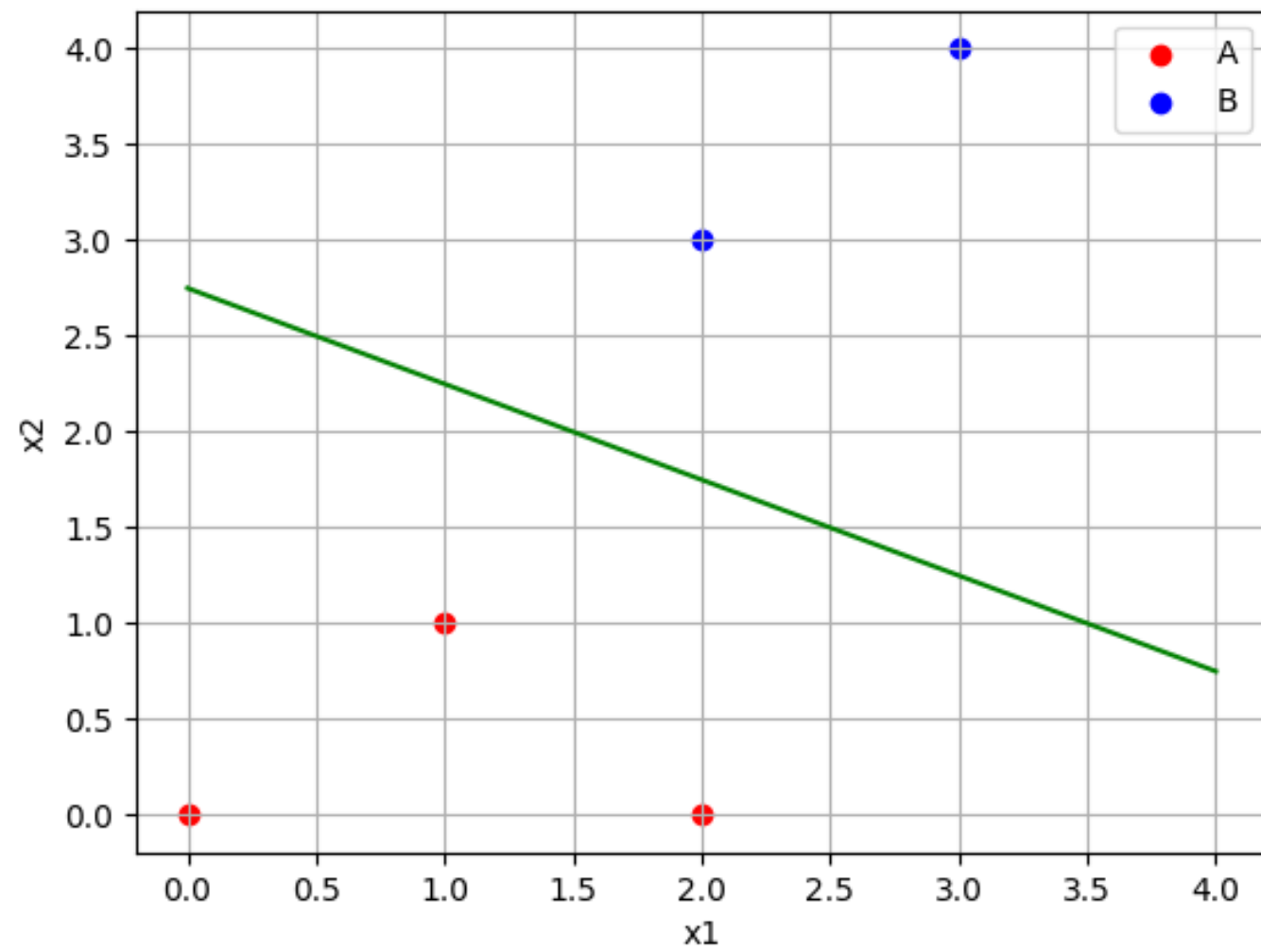
$$w_1x_1 + w_2x_2 + b = 0$$

$$w_1x_1 + w_2(0) + b = 0$$

$$w_1x_1 = -b$$

$$x_1 = -b / w_1$$





```

33 xx = np.linspace(0, 4)
34 yy = slope * xx - (b / w[1])
35 plt.plot(xx, yy, color='g')
36 plt.show()
37
38 print(model.predict([[3, 3]])) # B
39 print(model.predict([[4, 0]])) # A
40 print(model.predict([[2, 2]])) # B
41 print(model.predict([[1, 2]])) # A
42 print(model.predict([[1.5, 1.5]])) # A
43

```

Run: in_class_coding

C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"

Weights w1 and w2: [[0.4 0.8]]

Bias: [-2.2]

['B']

['A']

['B']

['A']

['A']

Process finished with exit code 0

Kernels

Transform input data into the required form: linear, polynomial, and radial basis function (RBF). Nonlinear kernels used to separate categories with curved hyperplanes.

Regularization

Parameter C controls the tolerance of classification error. A high C will seek to classify all points correctly, narrowing the margin. A low C will aim for the widest margin possible but will result in some points being misclassified.

Demo IRIS

Do the following:

Import the Iris data into a DataFrame. Print the columns, data types and number of non-null rows. Perform any data wrangling that is required.

FileEditViewNavigateCodeRefactorRunToolsVCSWindowHelp

ITP449_Fall2020 - in_class_coding.py

ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py

in_class_coding

in_class_coding.py

1import pandas as pd

2import os

3

4os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')

5pd.set_option('display.max_columns', None)

6

7irisDf = pd.read_csv('iris.csv')

8print(irisDf.info())

9

4271

Run: in_class_coding

Run

↑

↓

Undo

Redo

Copy

Paste

Close

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 150 entries, 0 to 149

Data columns (total 5 columns):

#	Column	Non-Null Count	Dtype
0	Sepal length	150 non-null	float64
1	Sepal width	150 non-null	float64
2	Petal length	150 non-null	float64
3	Petal width	150 non-null	float64
4	Species	150 non-null	object

dtypes: float64(4), object(1)

memory usage: 5.3+ KB

4: Run

TODO

6: Problems

Terminal

Python Console

Event Log

6:43 CRLF UTF-8 4 spaces Python 3.8 (Class)

Do the following:

What are the species categories?

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 pd.set_option('display.max_columns', None)
6
7 irisDf = pd.read_csv('iris.csv')
8 print(irisDf.info())
9 print(irisDf.Species.unique())
10
```

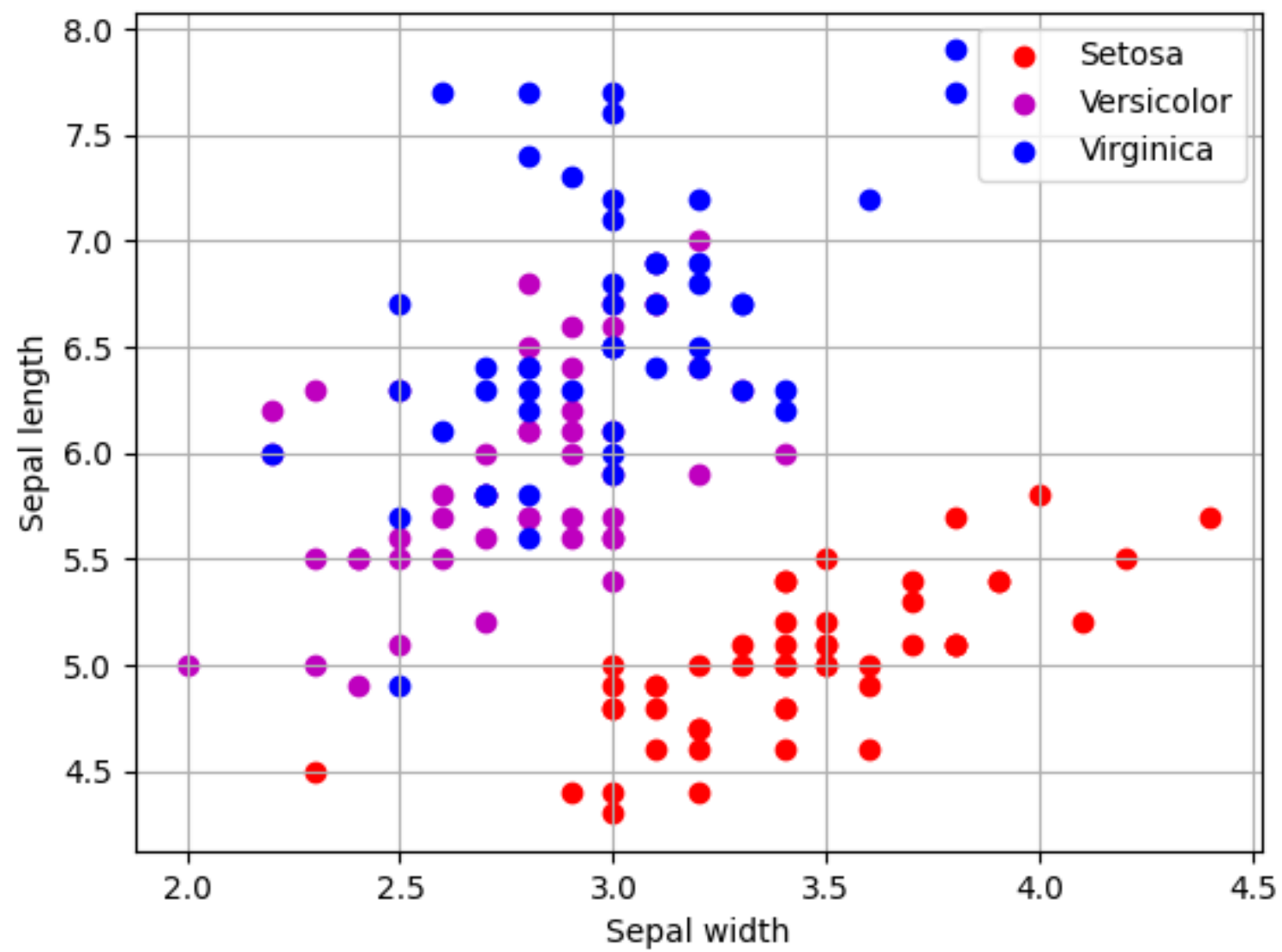
Run: in_class_coding

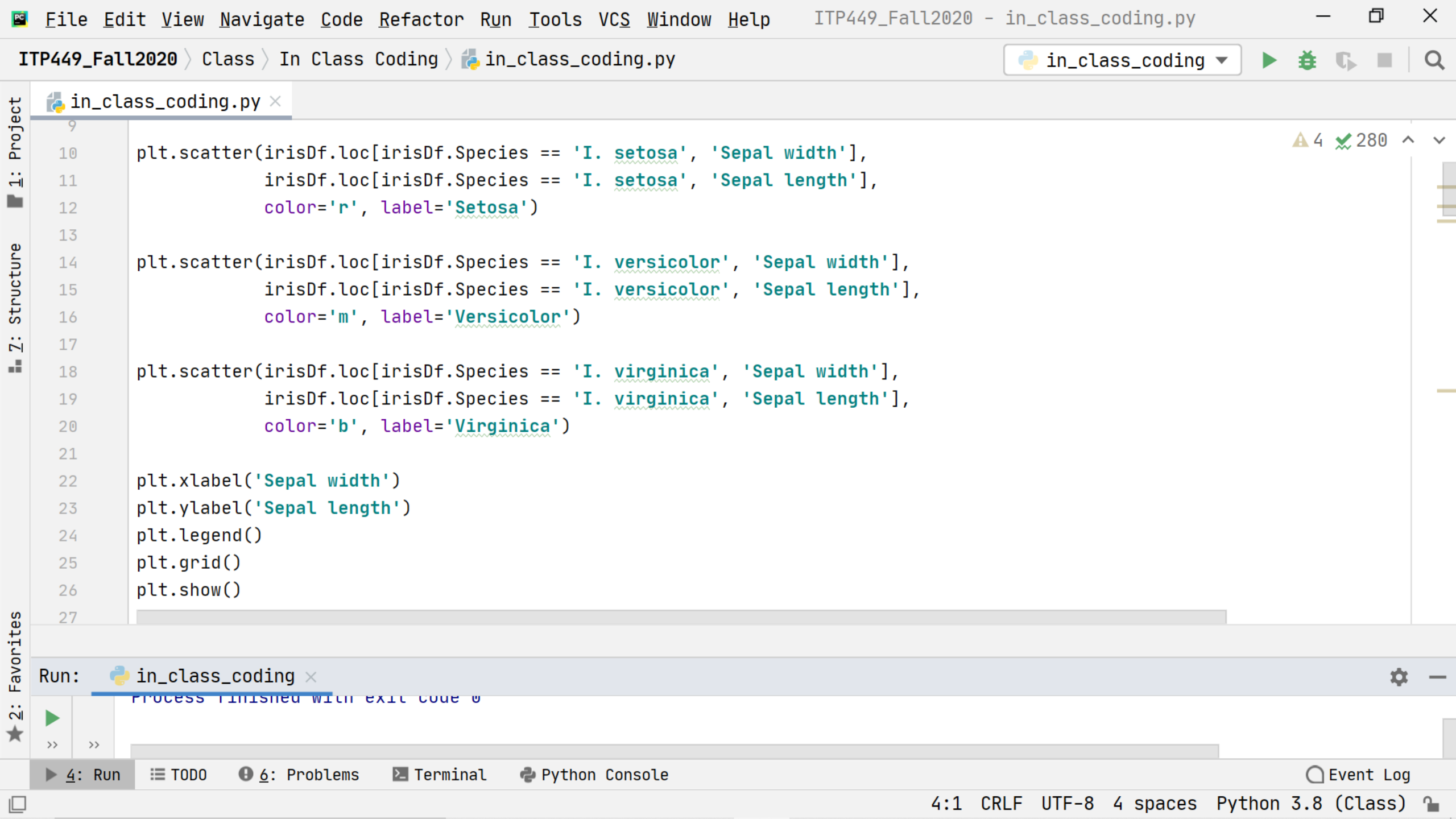
#	Column	Non-Null Count	Dtype
0	Sepal length	150 non-null	float64
1	Sepal width	150 non-null	float64
2	Petal length	150 non-null	float64
3	Petal width	150 non-null	float64
4	Species	150 non-null	object

dtypes: float64(4), object(1)
memory usage: 5.3+ KB
None
['I. setosa' 'I. versicolor' 'I. virginica']

Do the following:

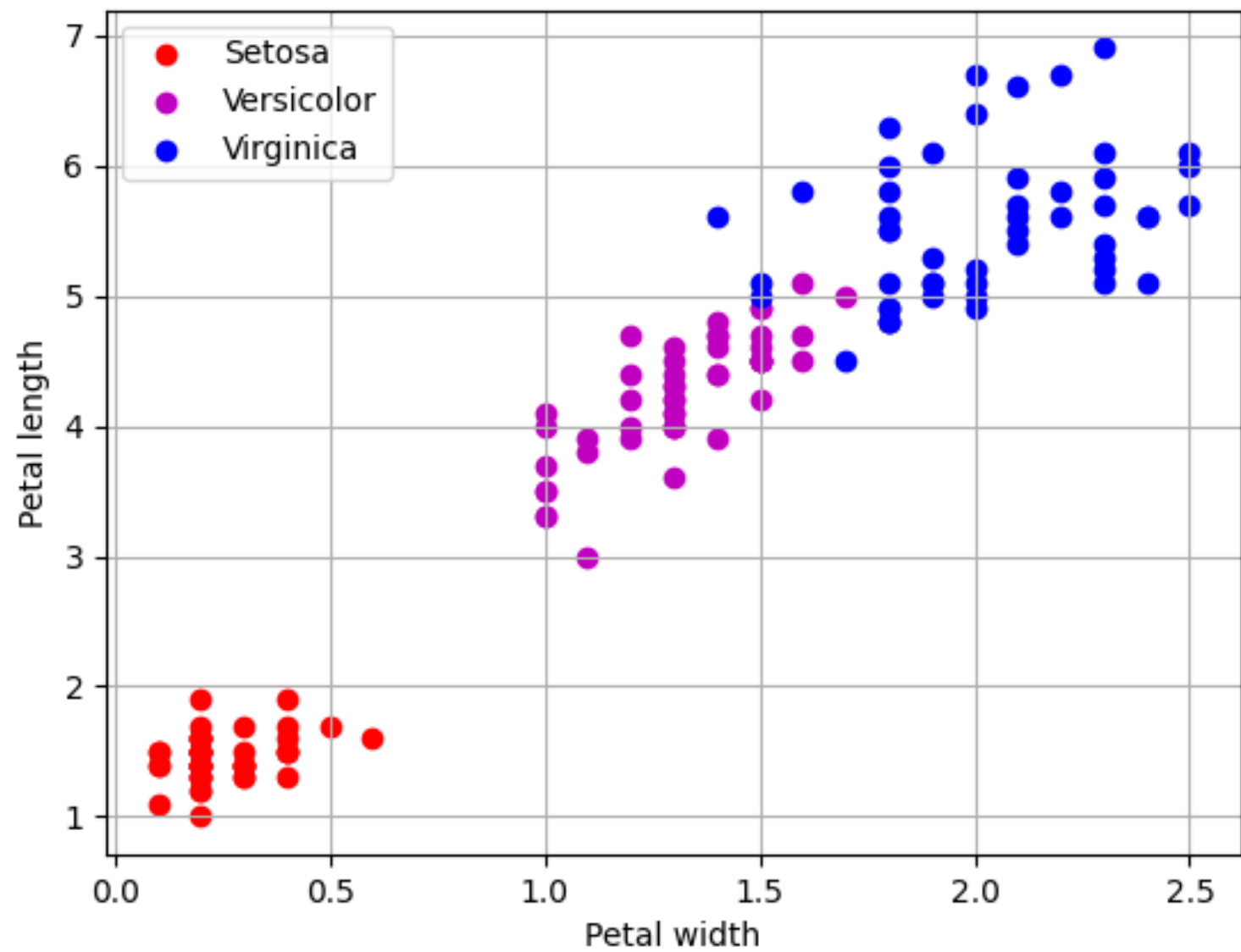
Plot a scatterplot of the sepal_width and sepal_length, assigning different colors to the species categories.
What is your observation?

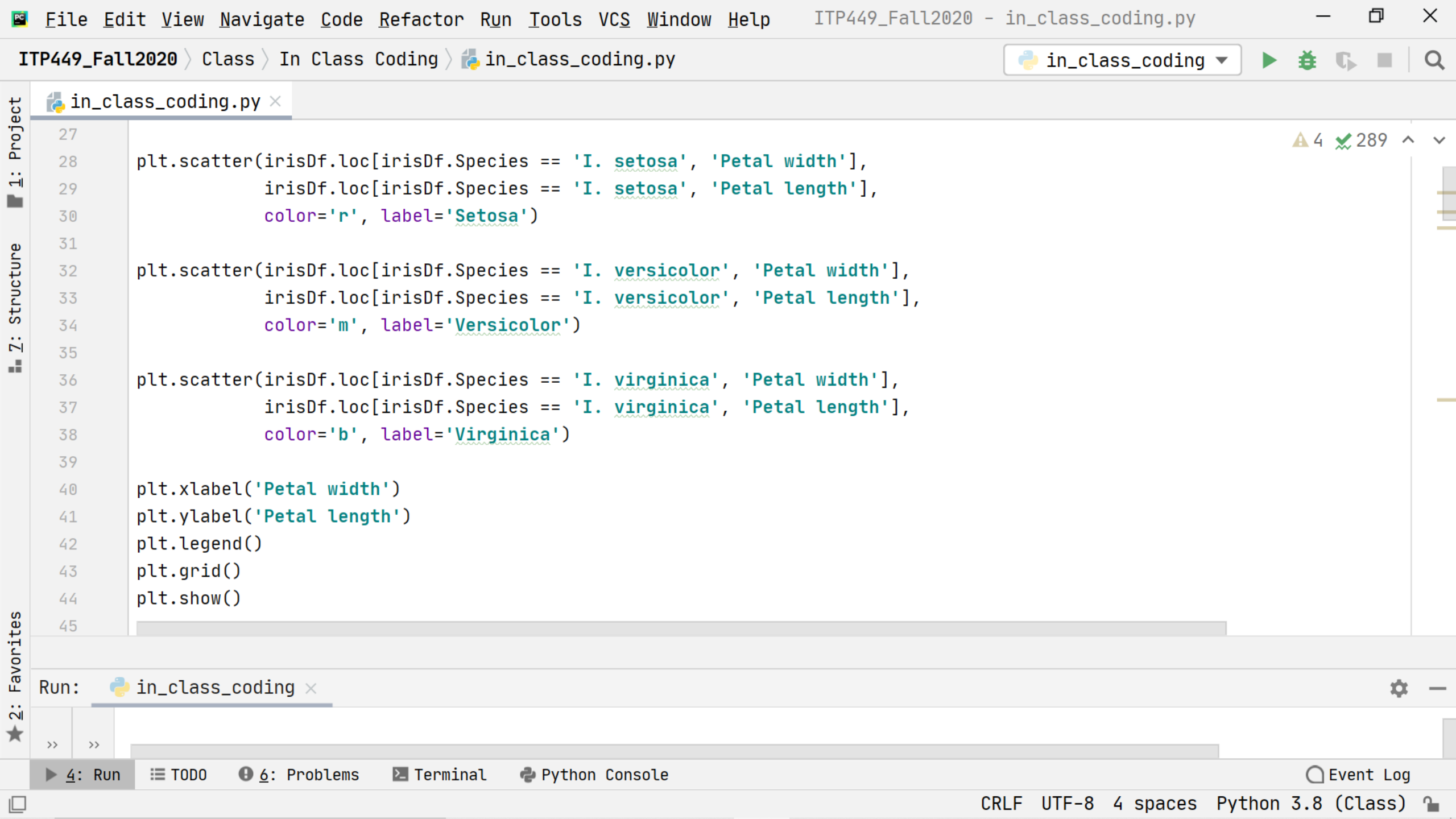




Do the following:

Plot a scatterplot of the `petal_width` and `petal_length`, assigning different colors to the species categories.
What is your observation?



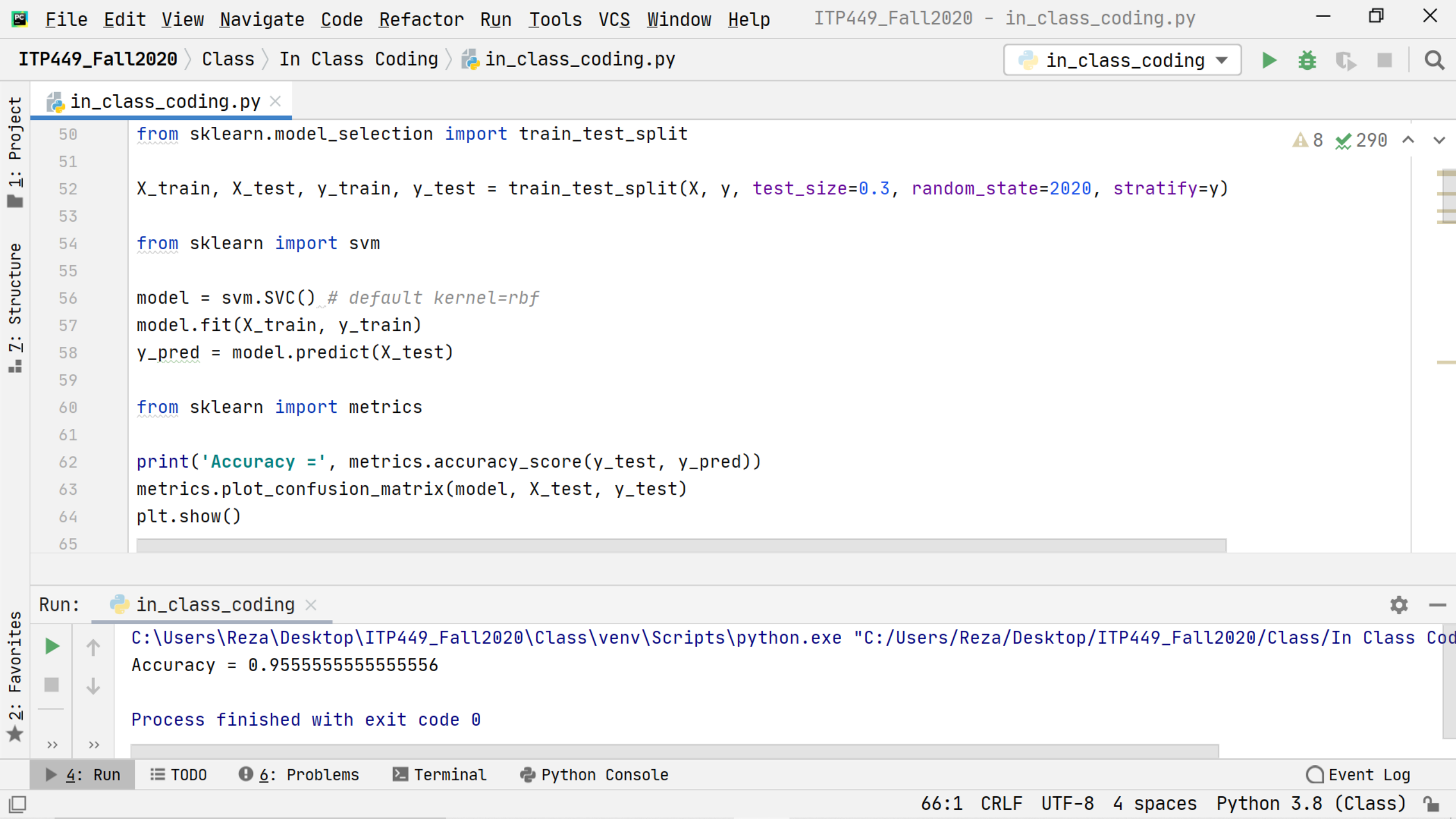


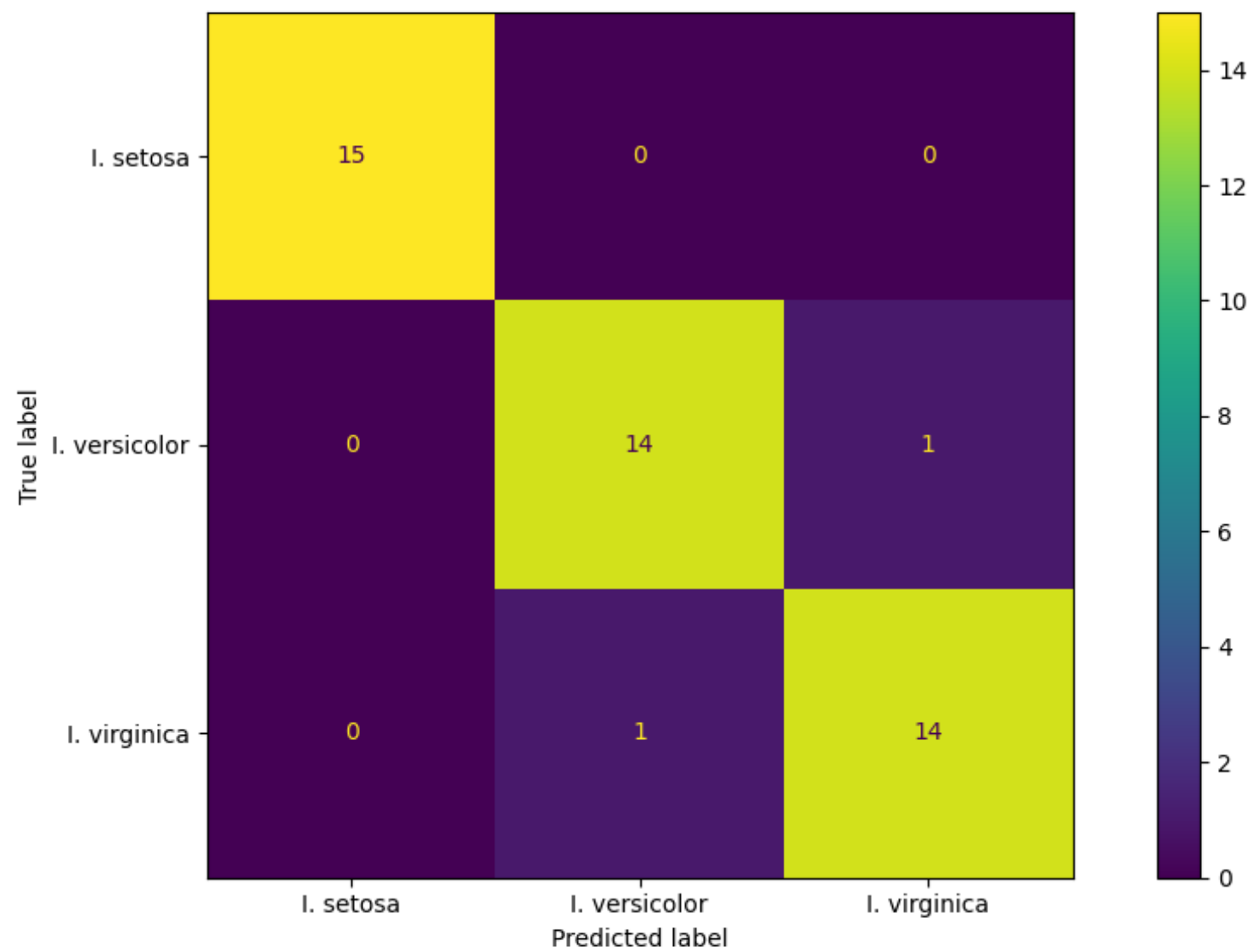
Do the following:

Split the data into training and testing sets. Train a model using all four features.

Plot the confusion matrix.

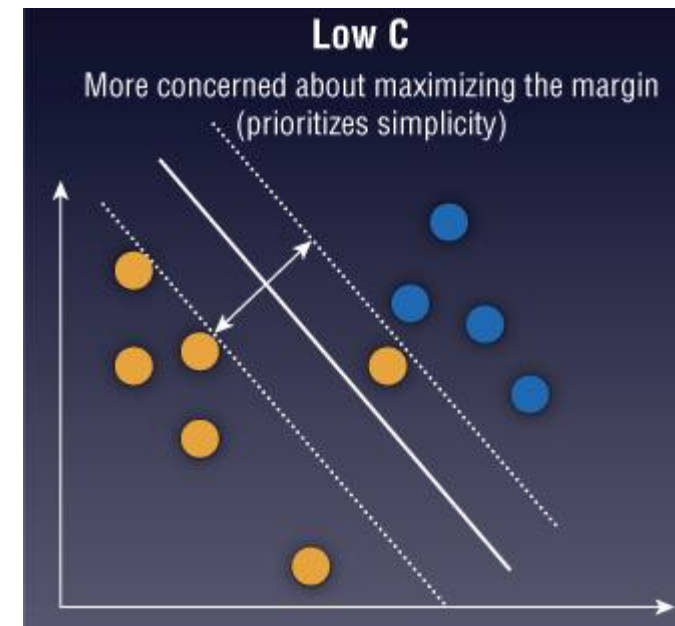
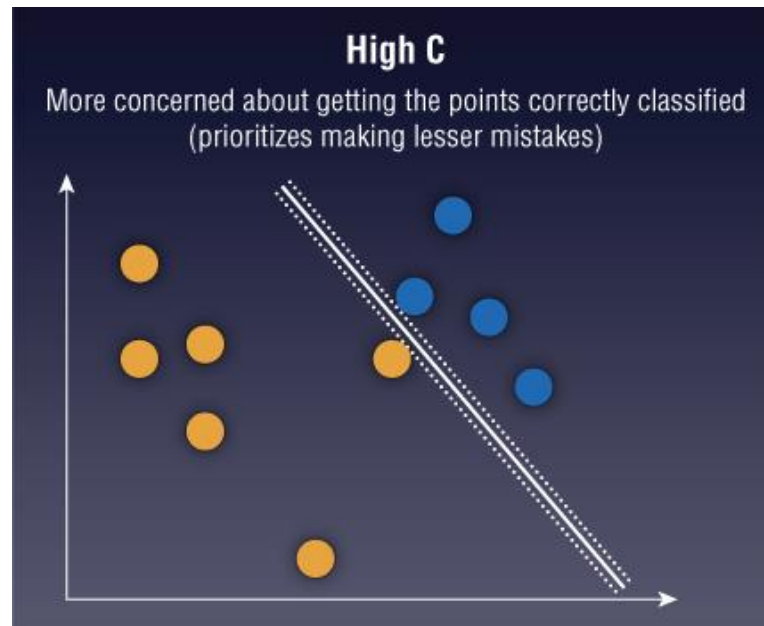
Print the accuracy score.





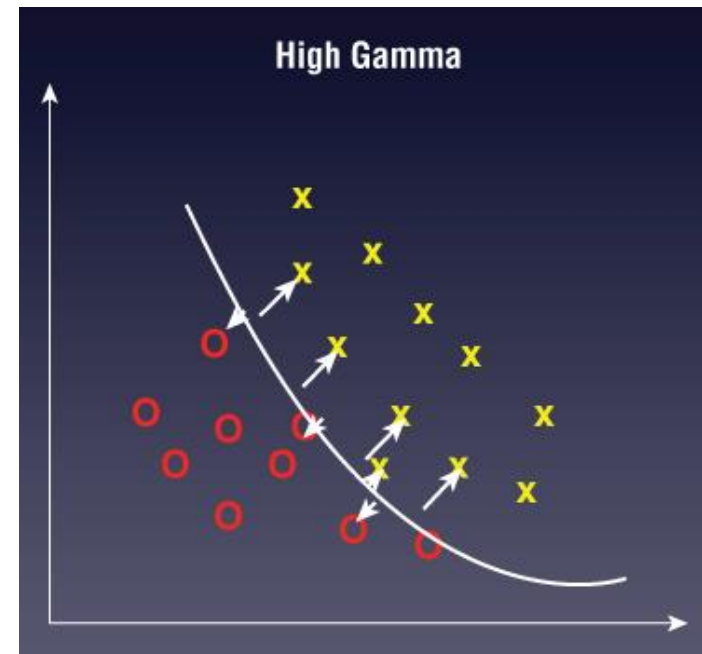
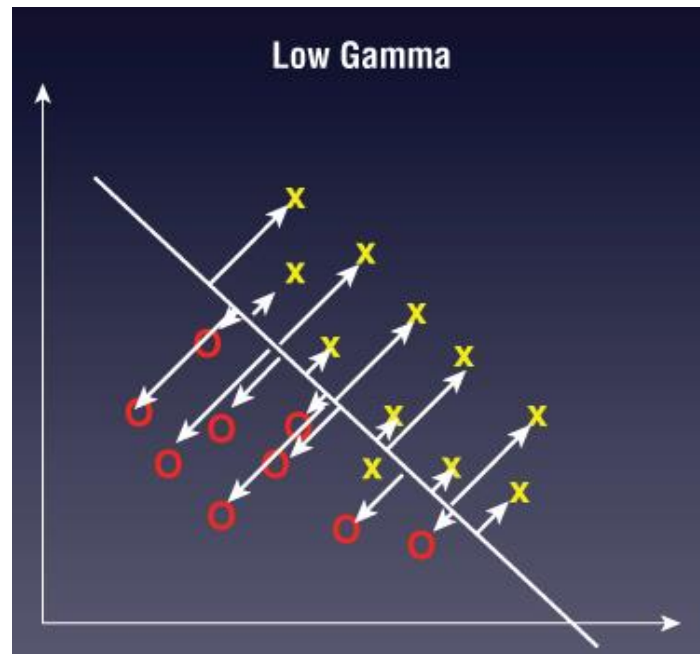
Regularization

Parameter C controls the tolerance of classification error. A high C will seek to classify all points correctly, narrowing the margin. A low C will aim for the widest margin possible but will result in some points being misclassified.



Gamma

Defines how far the influence of a single training example reaches. A low Gamma indicates that every point has a far reach while a high Gamma indicates that points closest to the decision boundary have a close reach.



Linear Kernel

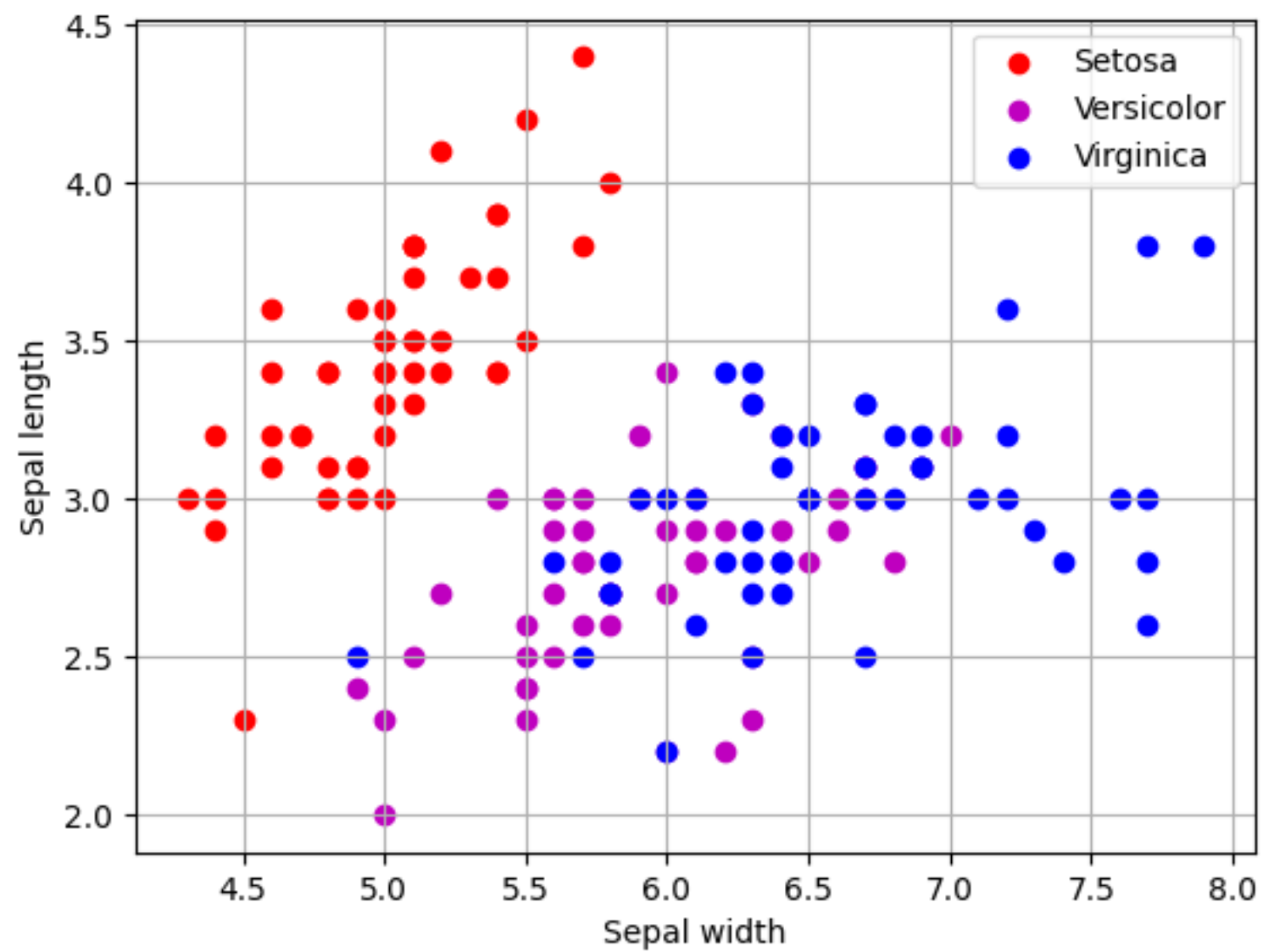
Uses a straight line to separate points for classification.

Polynomial Kernel

Uses polynomial curves to separate points for classification where a polynomial of degree 1 is essentially a linear kernel.

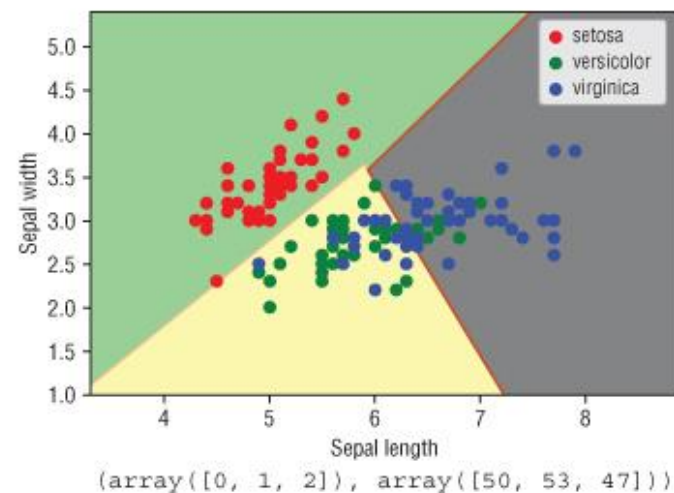
Radial Basis Function (RBF) Kernel

Gives value to each point based on its distance from the origin or a fixed center, commonly on a Euclidean space.



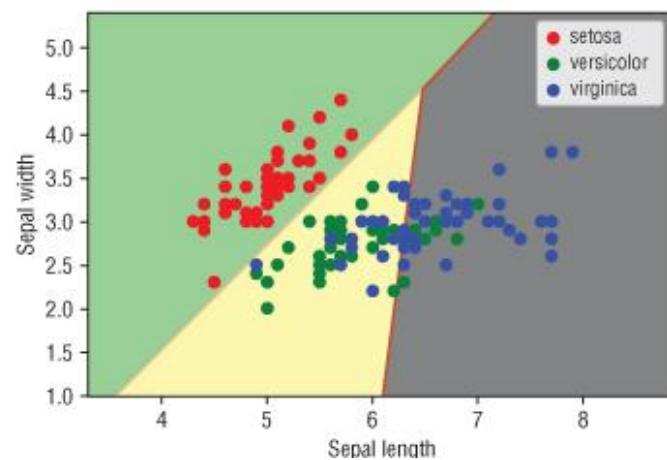
Default C C=1

SVC with linear kernel



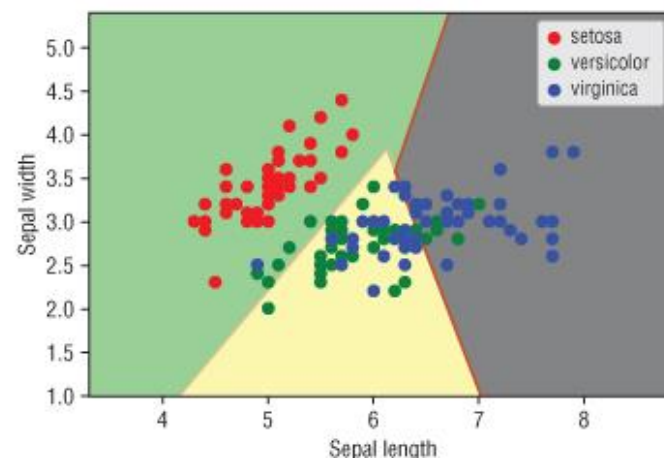
High C C=10¹⁰

SVC with linear kernel

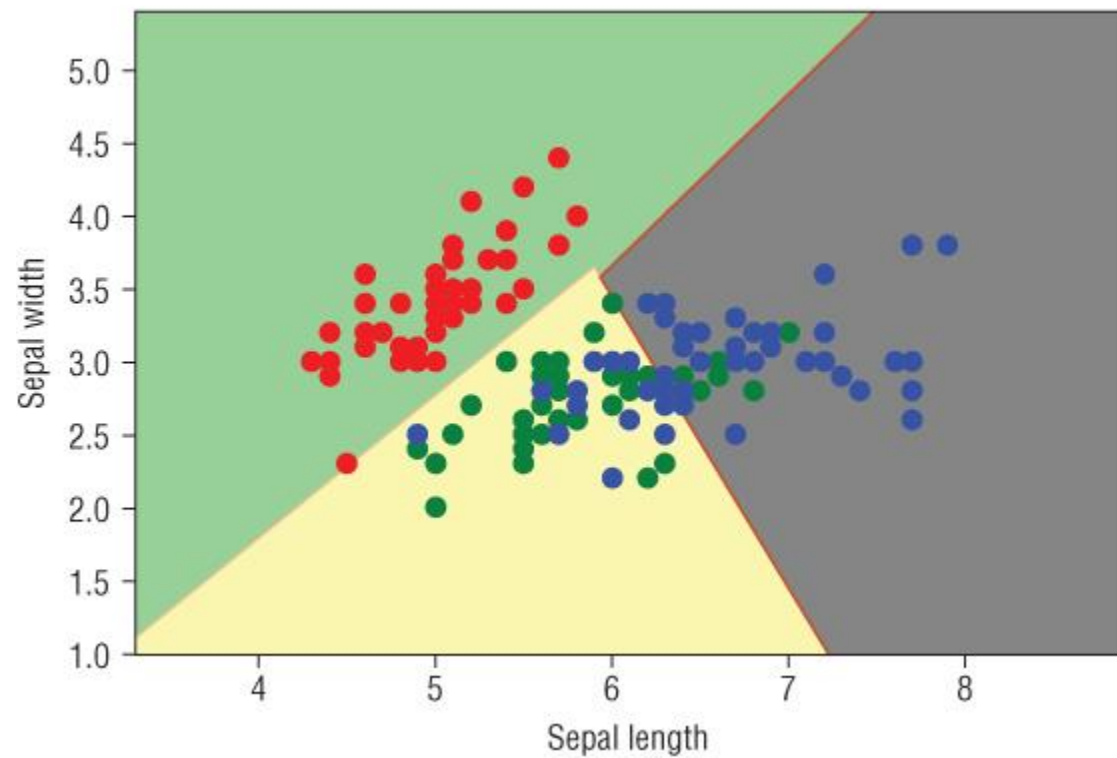


Low C C=10⁻¹⁰

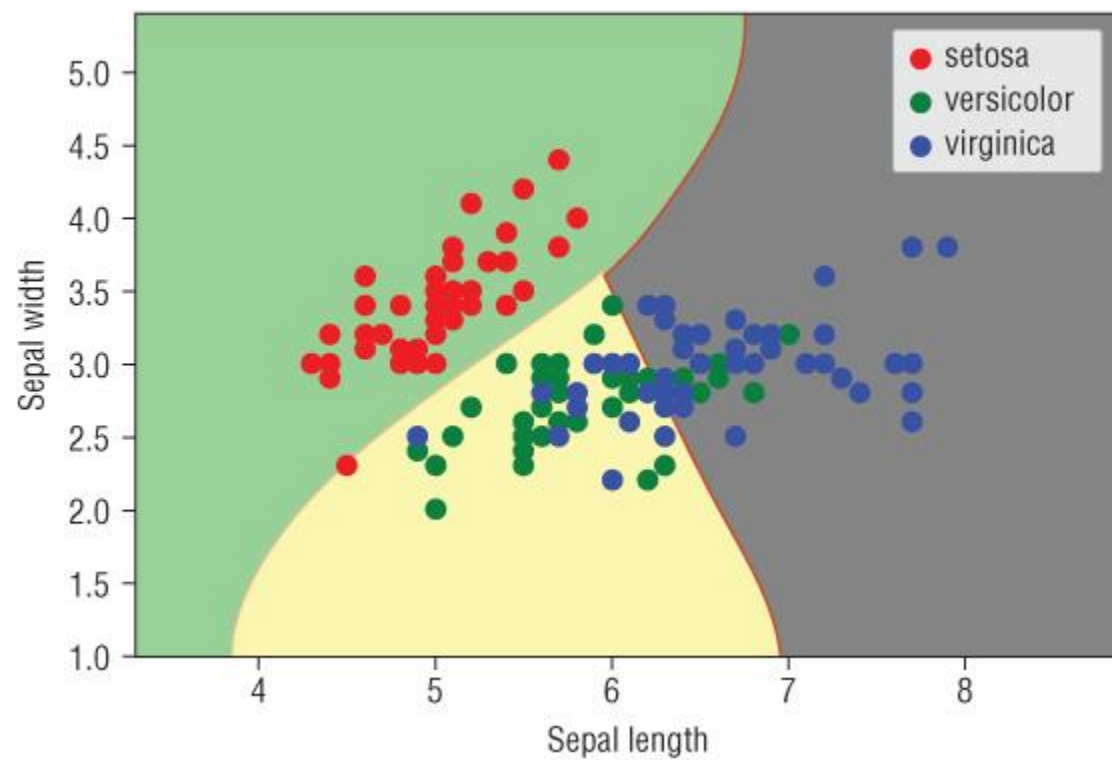
SVC with linear kernel



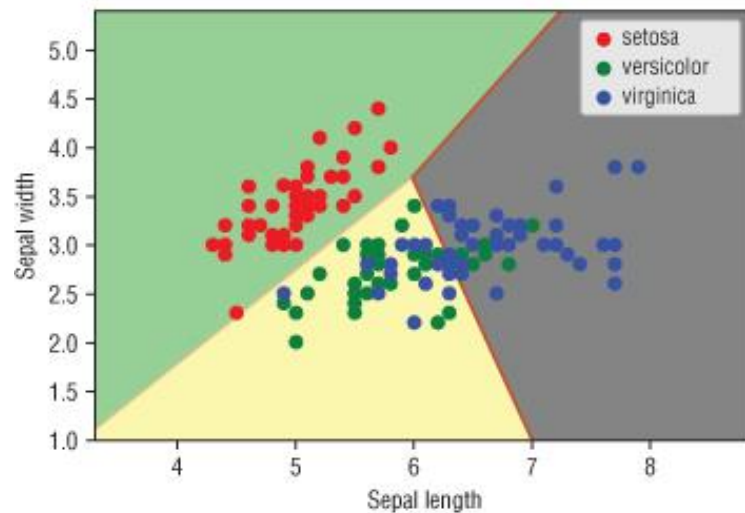
SVC with linear kernel



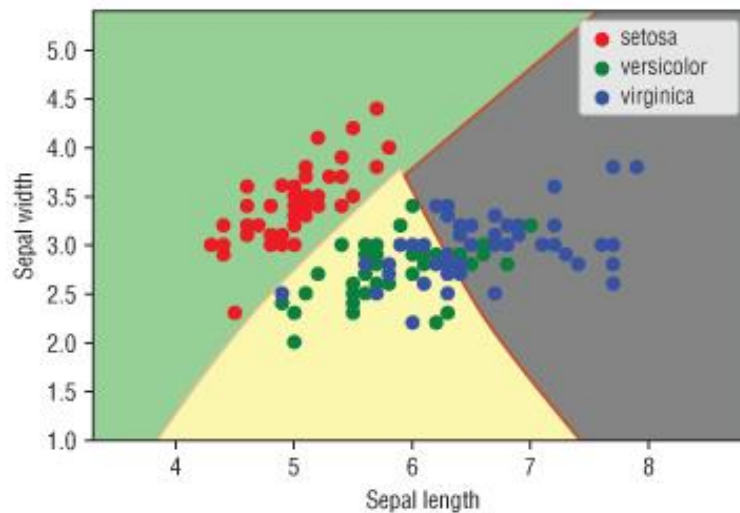
SVC with RBF kernel



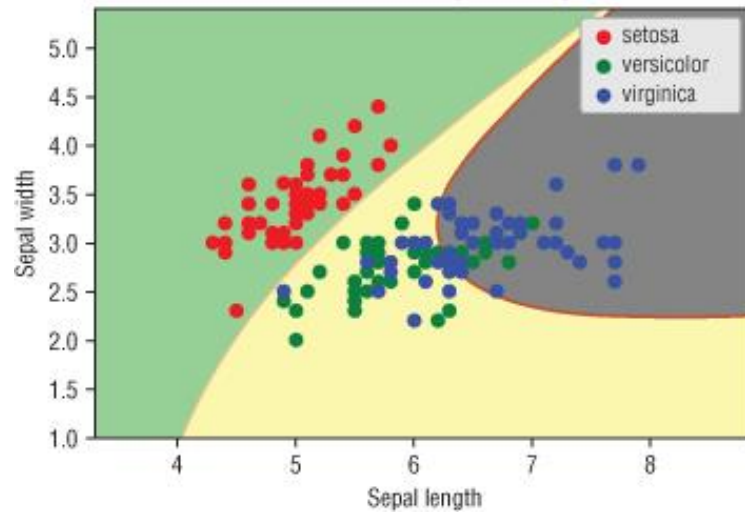
SVC with polynomial (degree 1) kernel



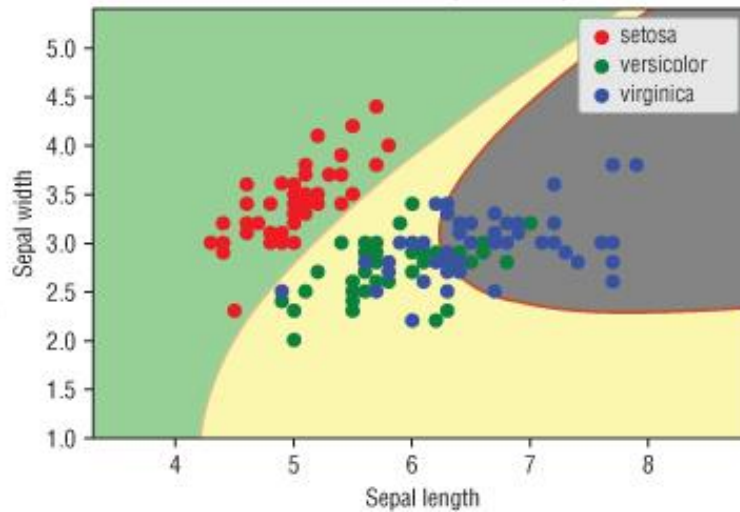
SVC with polynomial (degree 2) kernel



SVC with polynomial (degree 3) kernel



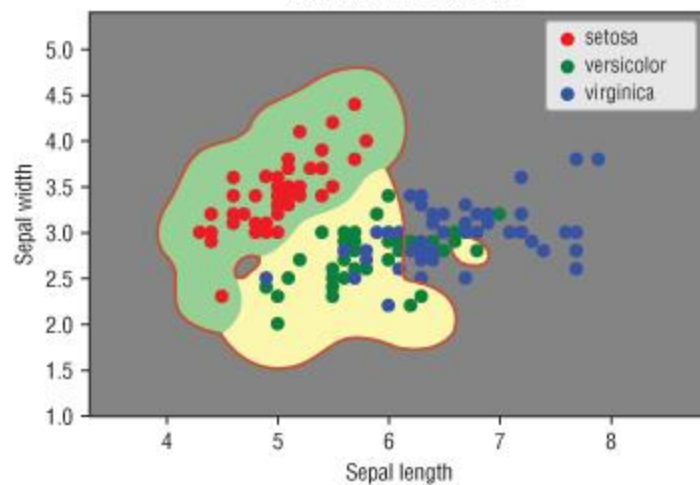
SVC with polynomial (degree 4) kernel



Gamma = 10

C = 1

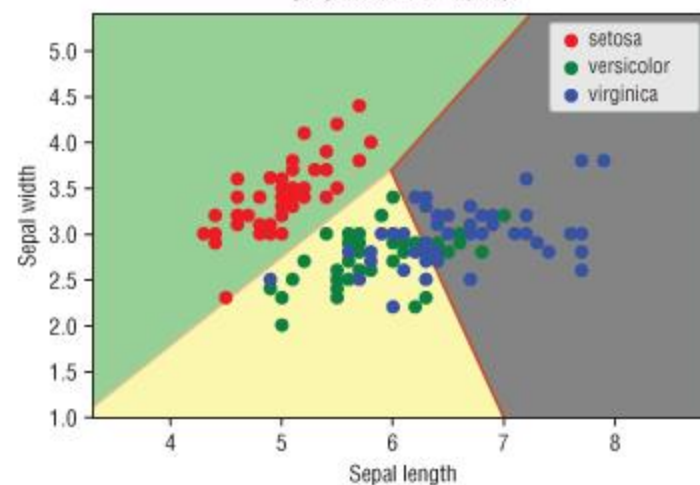
SVC with RBF kernel



Gamma = 0.1

C = 1

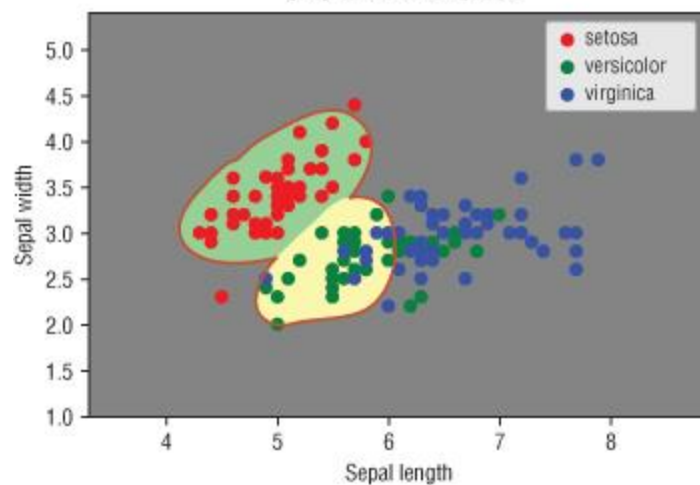
SVC with RBF kernel



Gamma = 10

C = 10⁻¹⁰

SVC with RBF kernel



Gamma = 0.1

C = 10¹⁰

SVC with RBF kernel

