

## Module 3-2 Exercises

The [Iris data set](#) is an extremely popular data set for learning machine learning techniques. The data set contains 3 classes that each correspond to a type of Iris plant. Each class has 50 points and each data point consists of 4 features.

These features consist of the sepal length in cm, sepal width in cm, petal length in cm, and petal width in cm.

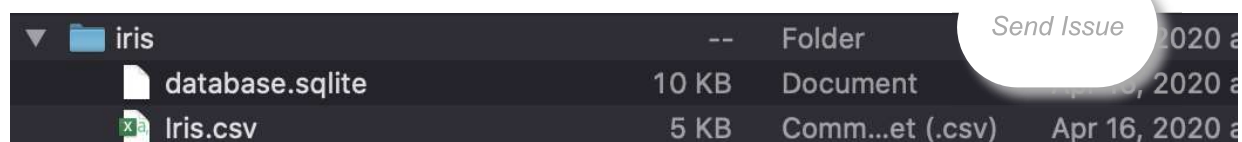
Today we are going to use svm from the machine learning library, sklearn, in order to predict the classes each Iris data point belongs to. Our goal is to pass the sepal length in cm, sepal width in cm, petal length in cm, and petal width in cm in our svm model and have it predict one of the three classes of Iris flower this data could belong to.

### Step 1

Before we can get started training our machine learning model we have to first get the data. Download the Iris data set.

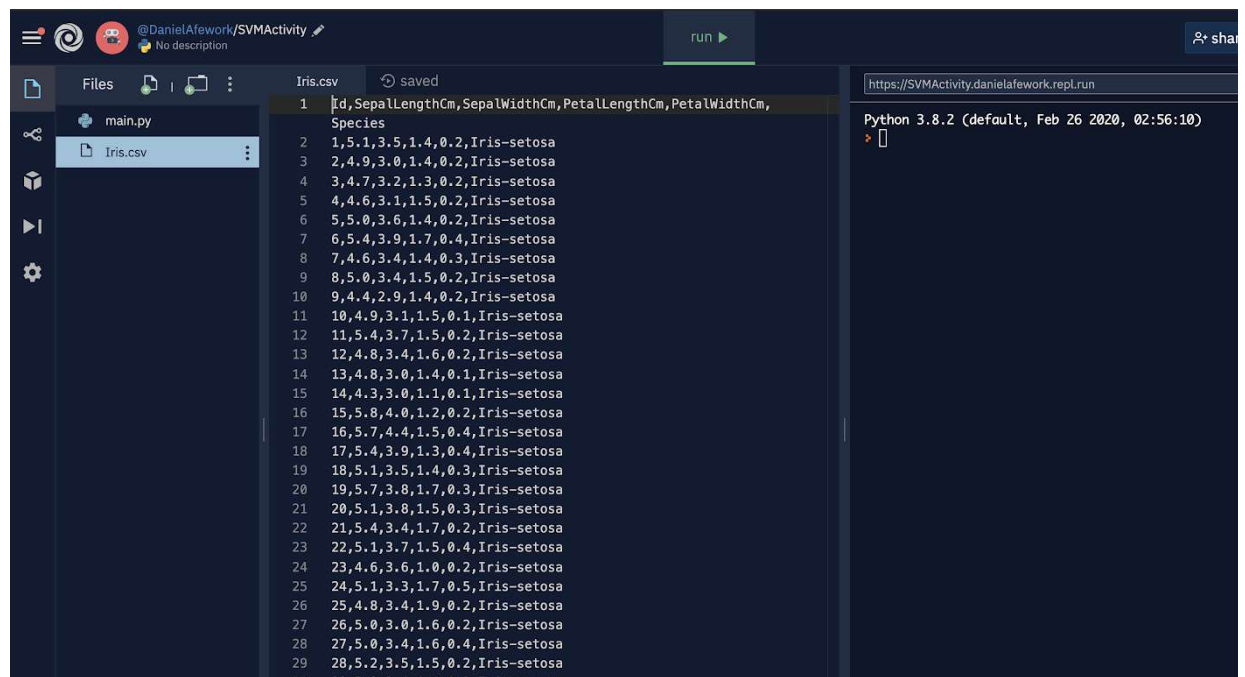


After clicking on the download button, a file called `archive.zip` should be downloaded to your computer. Unzip the file to reveal a folder containing two files, `database.sqlite` and `Iris.csv`.



From here, drag and drop `Iris.csv` into a new repl.it. If you are using a local installation of Python, drop it into whatever folder you will be making your program in. You should have

title "Iris.csv" inside of your repl which has contents that look like the following.



The screenshot shows a Repl.it workspace for a user named @DanielAFework. The workspace contains two files: `main.py` and `Iris.csv`. The `Iris.csv` file is open, displaying the following data:

```
1 Id,SepalLengthCm,SepalWidthCm,PetalLengthCm,PetalWidthCm,Species
2 1,5.1,3.5,1.4,0.2,Iris-setosa
3 2,4.9,3.0,1.4,0.2,Iris-setosa
4 3,4.7,3.2,1.3,0.2,Iris-setosa
5 4,4.6,3.1,1.5,0.2,Iris-setosa
6 5,5.0,3.6,1.4,0.2,Iris-setosa
7 6,5.4,3.9,1.7,0.4,Iris-setosa
8 7,4.6,3.4,1.4,0.3,Iris-setosa
9 8,5.0,3.4,1.5,0.2,Iris-setosa
10 9,4.4,2.9,1.4,0.2,Iris-setosa
11 10,4.9,3.1,1.5,0.1,Iris-setosa
12 11,5.4,3.7,1.5,0.2,Iris-setosa
13 12,4.8,3.4,1.6,0.2,Iris-setosa
14 13,4.8,3.0,1.4,0.1,Iris-setosa
15 14,4.3,3.0,1.1,0.1,Iris-setosa
16 15,5.8,4.0,1.2,0.2,Iris-setosa
17 16,5.7,4.4,1.5,0.4,Iris-setosa
18 17,5.4,3.9,1.3,0.4,Iris-setosa
19 18,5.1,3.5,1.4,0.3,Iris-setosa
20 19,5.7,3.8,1.7,0.3,Iris-setosa
21 20,5.1,3.8,1.5,0.3,Iris-setosa
22 21,5.4,3.4,1.7,0.2,Iris-setosa
23 22,5.1,3.7,1.5,0.4,Iris-setosa
24 23,4.6,3.6,1.0,0.2,Iris-setosa
25 24,5.1,3.3,1.7,0.5,Iris-setosa
26 25,4.8,3.4,1.9,0.2,Iris-setosa
27 26,5.0,3.0,1.6,0.2,Iris-setosa
28 27,5.0,3.4,1.6,0.4,Iris-setosa
29 28,5.2,3.5,1.5,0.2,Iris-setosa
30 29,5.2,3.4,1.4,0.2,Iris-setosa
```

## Step 2

Let's store the dataset inside of a pandas dataframe so that we can use it to train our model. The first step is going to be to first import a few helpful libraries inside of the `main.py` file.

Begin by adding these lines to the top of `main.py`.

Python

```
from sklearn import svm
```

```
import pandas as pd
```

The first line will add the `svm` library from `sklearn` so that we can use the `svm` classification model to train our data. The second line will add the `pandas` library so that we can use dataframes. `pandas` is a library used by data scientists that allows for easy data manipulation.

Now let's take the data from the `Iris.csv` file and add it to a `pandas` dataframe.

Python

```
df = pd.read_csv("Iris.csv")
```

```
print(df)
```

We can print the data frame by writing `print(df)` and get the following output. As you see, the information from the CSV file is now nicely organized inside of the `pandas` dataframe.

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...	...
145	146	6.7	3.0	5.2	2.3	Iris-versicolour
146	147	6.3	2.5	5.0	1.9	Iris-versicolour
147	148	6.5	3.0	5.2	2.0	Iris-versicolour
148	149	6.2	3.4	5.4	2.3	Iris-versicolour
149	150	5.9	3.0	5.1	1.8	Iris-versicolour

## Step 2

Let's take the "SepalLength", "SepalWidth", "PetalLength", and "PetalWidth" and separate them into a smaller data frame so that we can use them to train our model.

Python

```
X = df[['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm']]
```

Now let's do the same thing, but for the species of flower.

Python

```
Y = df[['Species']]
```

Print both X and Y out to confirm that their contents are the separated flower data and species data.

Python

```
print(X)
```

```
print(Y)
```

The output should look something like this.

```
      SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm
0              5.1             3.5           1.4           0.2
1              4.9             3.0           1.4           0.2
2              4.7             3.2           1.3           0.2
3              4.6             3.1           1.5           0.2
4              5.0             3.6           1.4           0.2
..              ...             ...           ...           ...
145             6.7             3.0           5.2           2.3
146             6.3             2.5           5.0           1.9
147             6.5             3.0           5.2           2.0
148             6.2             3.4           5.4           2.3
149             5.9             3.0           5.1           1.8

[150 rows x 4 columns]
      Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
..      ...
145  Iris-virginica
146  Iris-virginica
147  Iris-virginica
148  Iris-virginica
149  Iris-virginica

[150 rows x 1 columns]
```

## Step 3

Now, before we train our model, we have to encode the species data into numerical values. The species data is currently stored as a `string` value. However, the machine learning model expects the data to be a numerical value.

To convert, we will use a `LabelEncoder` . Add this line of code to the top of your program

Python

```
from sklearn.preprocessing import LabelEncoder
```

Now we will encode the data using the `LabelEncoder` .

Python

```
le = LabelEncoder()

yEncoded = le.fit_transform(Y['Species'])
```

Note that we passed in `"Species"` instead of just `Y` . This is because the `le.fit_transform` function is expecting a list and not the whole dataframe. `Y['Species']` will return the values underneath the species column.

If we print `yEncoded` , the following values will be returned.

```
[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
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2]
```

## Step 4

The original text data for species has been encoded as numerical data! Now let's train our model

Python

```
IrisPredictionModel = svm.SVC()

IrisPredictionModel.fit(X, yEncoded)
```

Now that our model is trained let's try to predict a flower.

Python

```
prediction = IrisPredictionModel.predict([ [4.9, 3, 1.4, .2]] )  
  
print(prediction)
```

We are trying to predict a flower with **SepalLengthCm** of 4.9 , **SepalWidthCm** of 3 , **PetalLengthCm** of 1.4 , and a **PetalWidthCm** of .2. After printing the prediction, we receive this value as



```
[0]
```

## Step 5

Now, our final step will be to translate the encoded number back to its corresponding value.

Python

```
returnToOriginal = le.inverse_transform(prediction)  
  
print(returnToOriginal)
```

This is the final value that will be printed out after returning the encoded data back to its original non numerical data.



```
['Iris-setosa']
```

Try predicting different numbers for the sepal length, width, and petal length and observe the results.