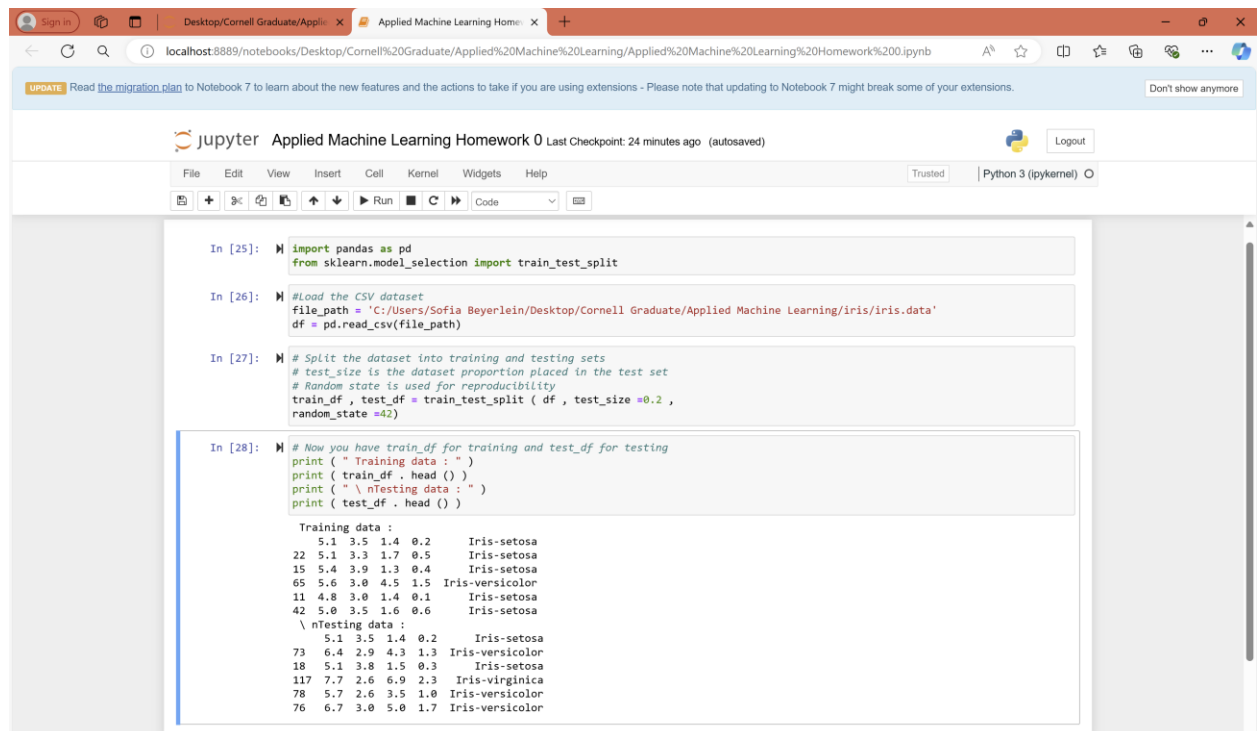# Applied Machine Learning Homework 0

1. There are 4 attributes features/attributes per sample. There are 3 different species (Iris Setosa, Iris Versicolour, and Iris Virginica) and there's 50 samples per each specie so a total of 150 instances.

2. Figure out how to parse the dataset you downloaded. Load the samples into an N × p array, where N is the number of samples and p is the number of attributes per sample. Additionally, create a N -dimensional vector containing each sample's label (species).



3. I plotted all the graphs using this snippet of code but replacing plt.scatter(sepal_length, sepal_width, color) with the other attribute arrays (e.g. plt.scatter(sepal_length, petal_length, color). I also want to clarify that I am

plotting the training data.
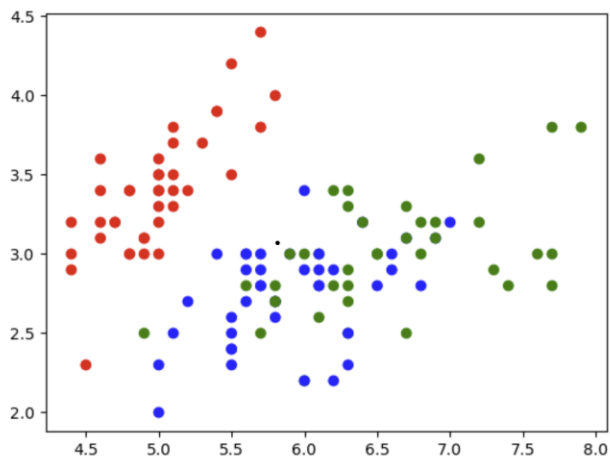
```
In [57]:  ▶  #saving the attributes of training data
             sepal_length = train_df.iloc[:, 0]
             sepal_width = train_df.iloc[:, 1]
             petal_length = train_df.iloc[:, 2]
             petal_width = train_df.iloc[:, 3]

             species = train_df.iloc[:, 4]        •
             color = []

             #assigning an rbg value to the species in the array
             for specie in species:
                 if specie == "Iris-setosa":
                     color.append("r")
                 if specie == "Iris-versicolor":
                     color.append("b")
                 if specie == "Iris-virginica":
                     color.append("g")

             plt.scatter(sepal_length, sepal_width, c=color)
             plt.savefig("plot.png")
```
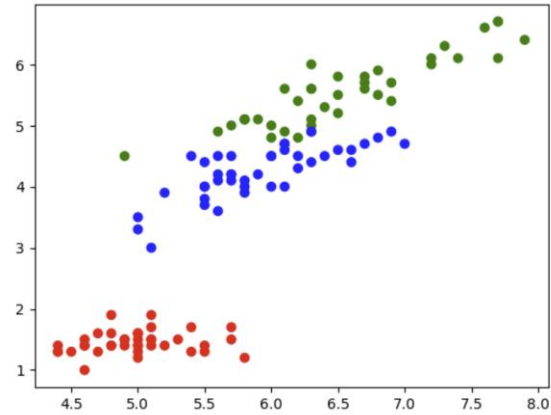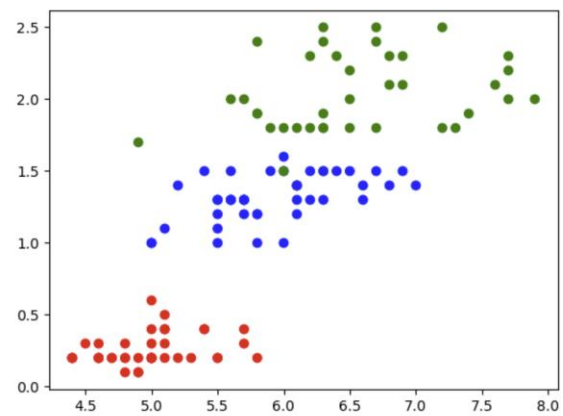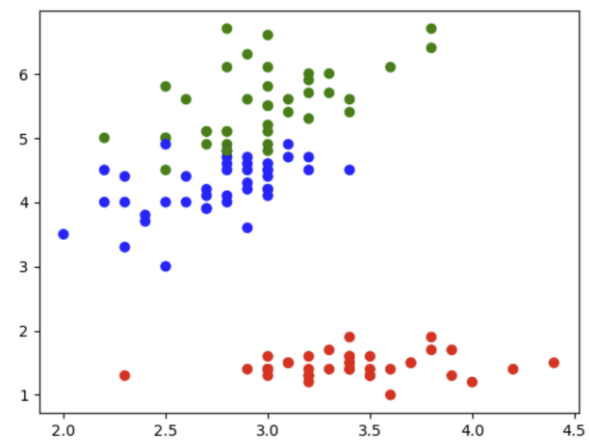
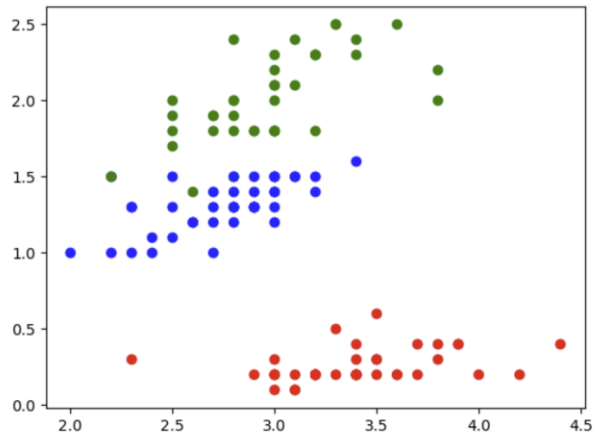Sepal Length vs. Sepal Width

## Sepal Length vs. Petal Length



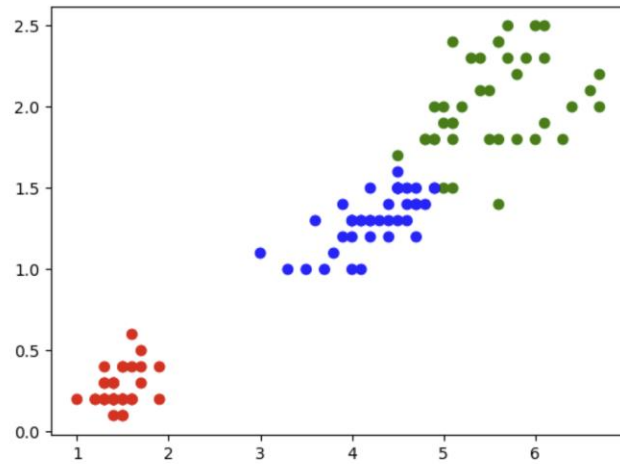## Sepal Length vs. Petal Width



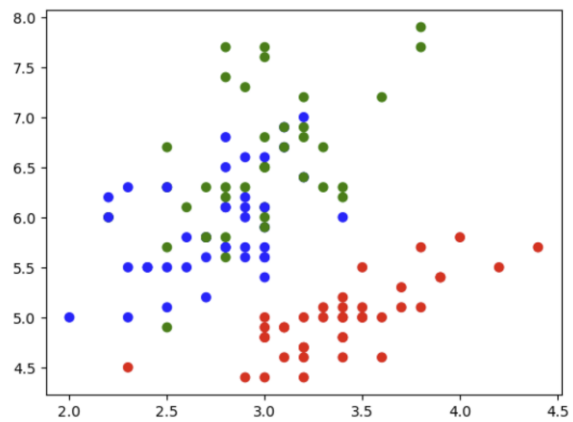## Sepal Width vs. Petal Length
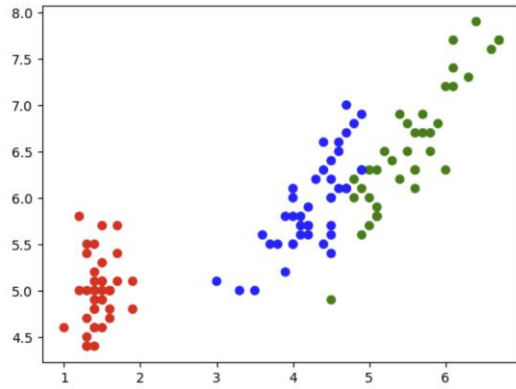
## Sepal Width vs. Petal Width
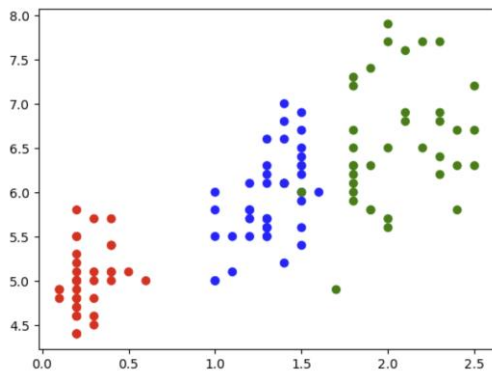


## Petal Length vs. Petal Width
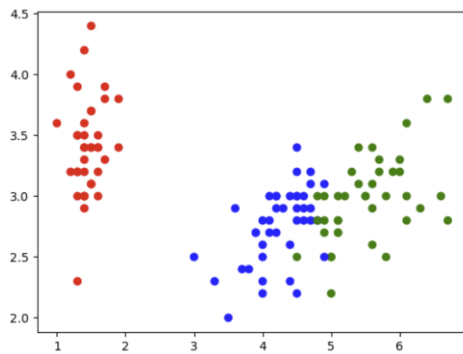


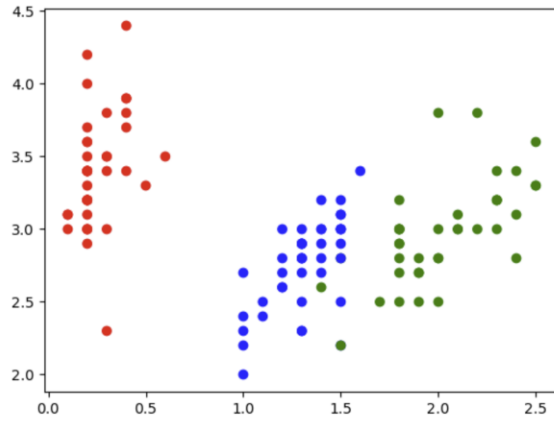## Sepal Width vs. Sepal Length



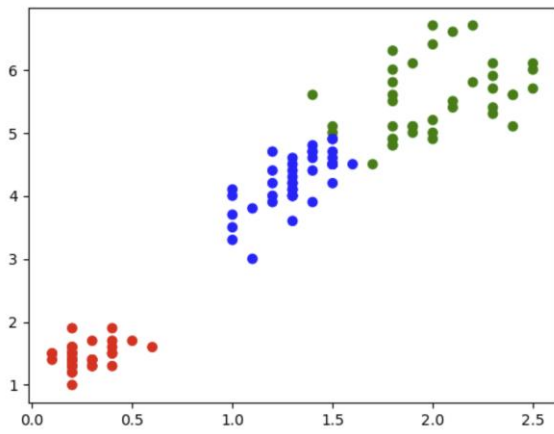## Petal Length vs. Sepal Length

Petal Width vs. Sepal Length
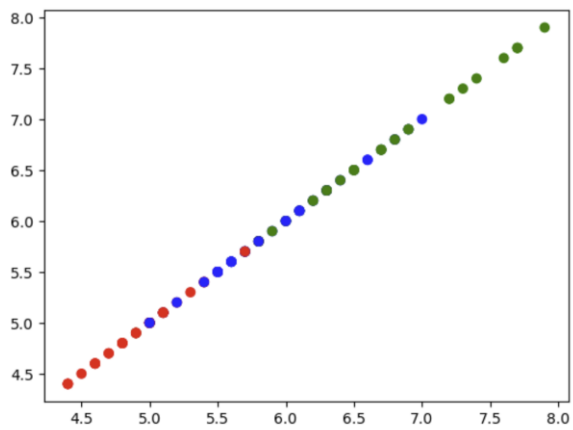


Petal Length vs. Sepal Width
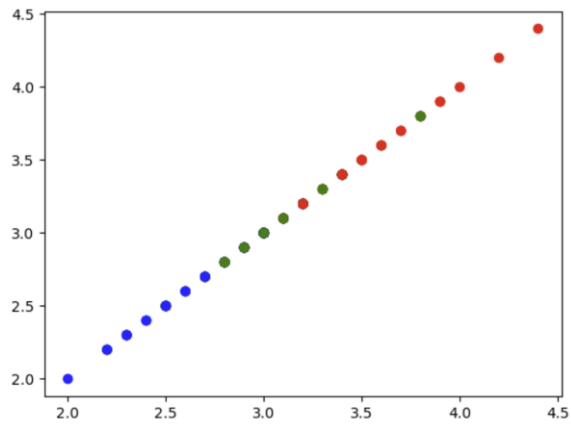


Petal Width vs. Sepal Width
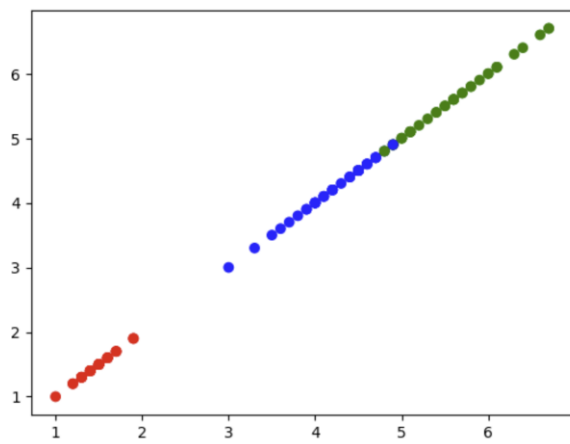
Petal Width vs. Petal Length



Sepal Length vs. Sepal Length



Sepal Width vs. Sepal Width

Petal Length vs. Petal Length



Petal Width vs. Petal Width