

Plotting for Exploratory data analysis (EDA)

Iris Flower dataset

Description

The Iris dataset was used in R.A. Fisher's classic 1936 paper. The Use of Multiple Measurements in Taxonomic Problems, and can also be found on the UCI Machine Learning Repository.

It includes three iris species with 50 samples each as well as some properties about each flower. One flower species is linearly separable from the other two, but the other two are not linearly separable from each other.

The columns in this dataset are:

Id Sepal.LengthCm Sepal.WidthCm Petal.LengthCm Petal.WidthCm Species

```
In [5]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

'''download iris.csv from https://raw.githubusercontent.com/uluc-cse/data-fair4/gh-pages/data/iris.csv
load iris.csv into a pandas dataframe.
iris = pd.read_csv("C:\Users\WAGARJUNA\Desktop\Downloads\Iris.csv")

In [6]: #Q1 How many data-points and features?
print (iris.shape)

(150, 5)

In [7]: #Q2 What are the column names in our dataset?
print (iris.columns)

Index(['sepal_length', 'sepal_width', 'petal_length', 'petal_width',
       'species'],
      dtype='object')

In [8]: #Q3 How many data points for each class are present?
#Q4 How many flowers for each species are present?
iris['species'].value_counts()
# balanced-dataset vs imbalanced dataset
# Iris is a balanced dataset as the number of data points for every class is 50.

Out[8]: virginica      50
         versicolor   50
         setosa       50
         Name: species, dtype: int64
```

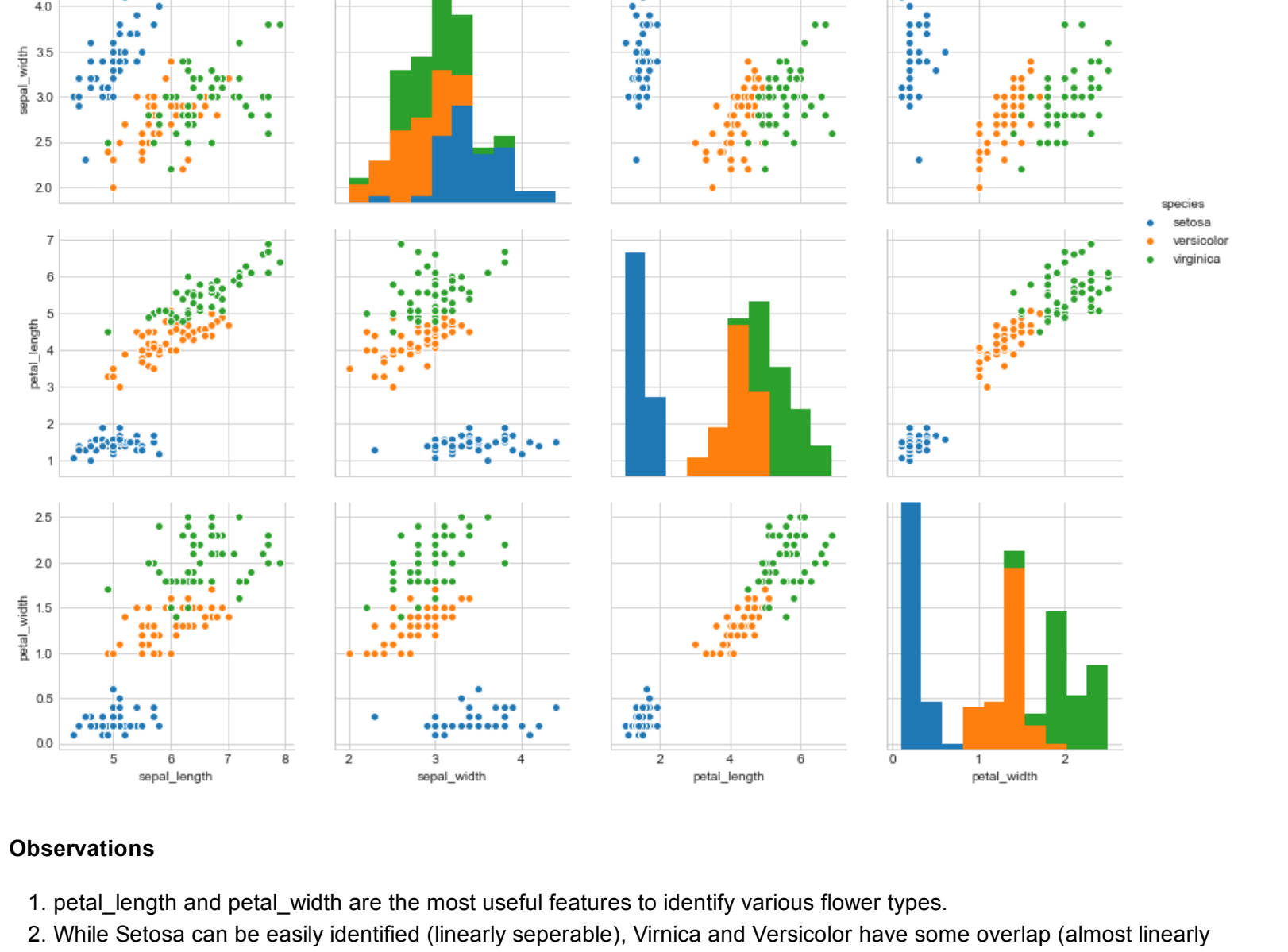
2-D Scatter Plot

```
In [9]: #2-D scatter plot:
#Always understand the axis: labels and scale.
iris.plot(kind='scatter', x='sepal_length', y='sepal_width') ;
plt.show()

# doesn't make much sense out of it.
# What if we color the points by their class-label/flower-type.

In [10]: # 2-D Scatter plot with color-coding for each flower type/class.
# Name 'sns' corresponds to seaborn.
sns.set_style('whitegrid')
sns.FacetGrid(iris, hue='species', size=4) \
    .map(plt.scatter, "sepal_length", "sepal_width") \
    .add_legend() \
    .plt.show()

# Notice that the blue points can be easily separated
# from red and green data points cannot be easily separated.
# Can we draw multiple 2-D scatter plots for each combination of features?
# How many combinations exist? 4C2 = 6.
```

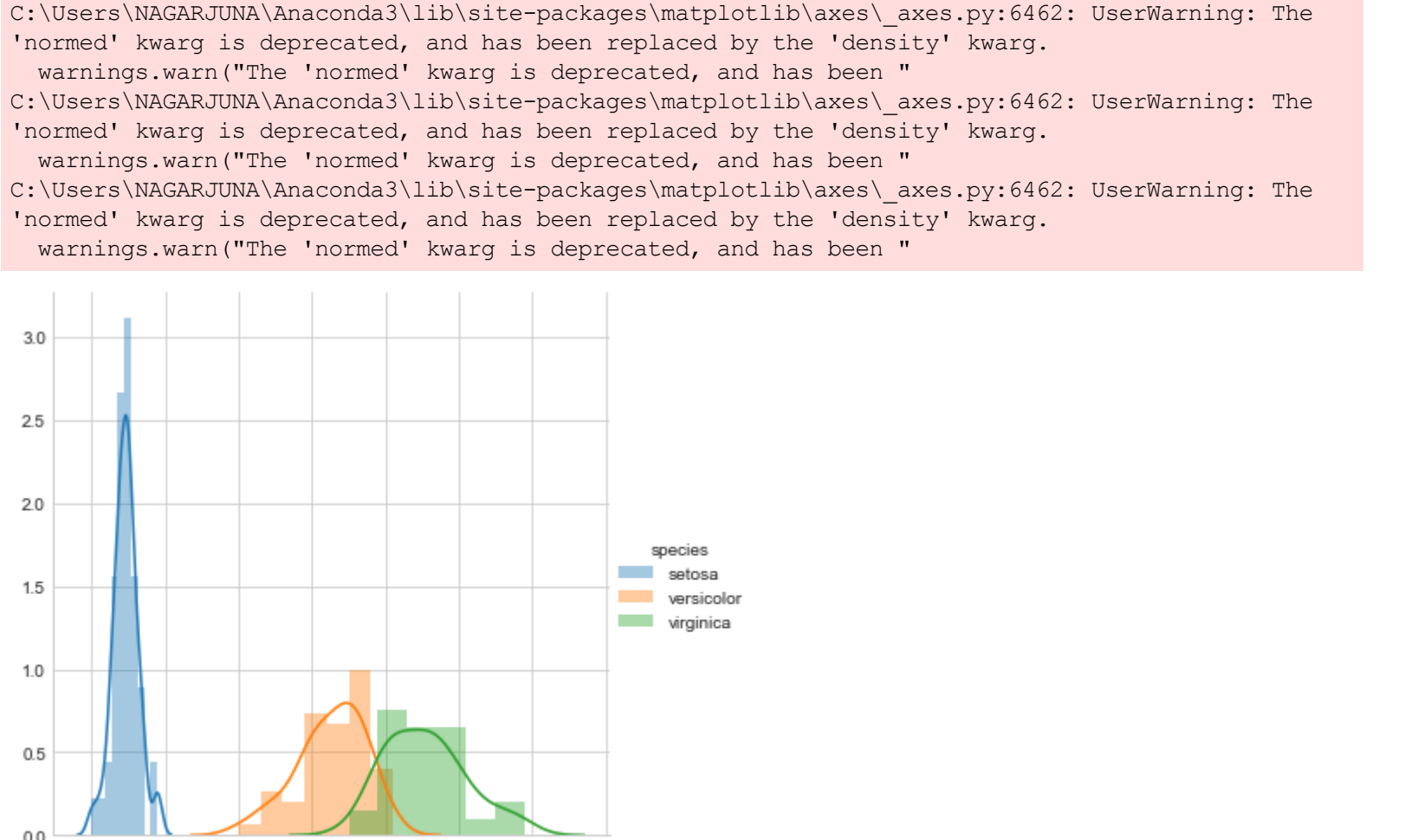


- Observations:**
1. Using sepal_length and sepal_width features, we can distinguish Setosa flowers from others.
 2. Separating Versicolor from Virginica is much harder as they have considerable overlap.
 3. We can find "lines" and "if-else" conditions to build a simple model to classify the flower types.

Pair-plot

```
In [11]: # pairwise scatter plot: Pair-Plot
# Dis-advantages:
# Can be used when number of features are high.
# Cannot visualize higher dimensional patterns in 3-D and 4-D.
# Only possible to view 2D patterns.
plt.close()
sns.set_style('whitegrid')
sns.pairplot(iris, hue='species', size=3)
plt.show()

# Notice: the diagonal elements are PDFs for each feature. PDFs are explained below.
```



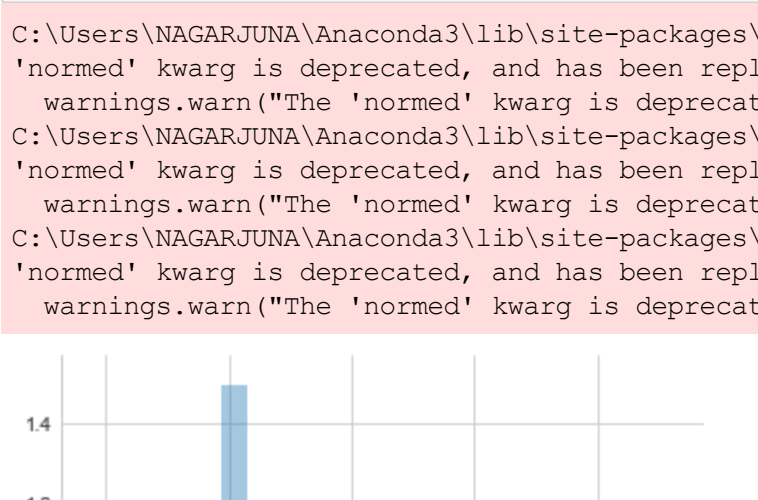
- Observations**
1. petal_length and petal_width are the most useful features to identify various flower types.
 2. While Setosa can be easily identified (linearly separable), Versicolor and Virginica have some overlap (almost linearly separable).
 3. We can find "lines" and "if-else" conditions to build a simple model to classify the flower types.

Histogram, PDF, CDF

```
In [12]: # What about 1-D scatter plot using just one feature?
# 1-D scatter plot of petal-length
import numpy as np
iris_setosa = iris.loc[iris['species'] == "setosa"]
iris_virginica = iris.loc[iris['species'] == "virginica"]
iris_versicolor = iris.loc[iris['species'] == "versicolor"]
# print(iris_setosa["petal_length"])
# print(iris_virginica["petal_length"])
# print(iris_versicolor["petal_length"])
plt.plot(iris_setosa["petal_length"], np.zeros_like(iris_setosa["petal_length"]), 'o')
plt.plot(iris_versicolor["petal_length"], np.zeros_like(iris_versicolor["petal_length"]), 'o')
plt.plot(iris_virginica["petal_length"], np.zeros_like(iris_virginica["petal_length"]), 'o')

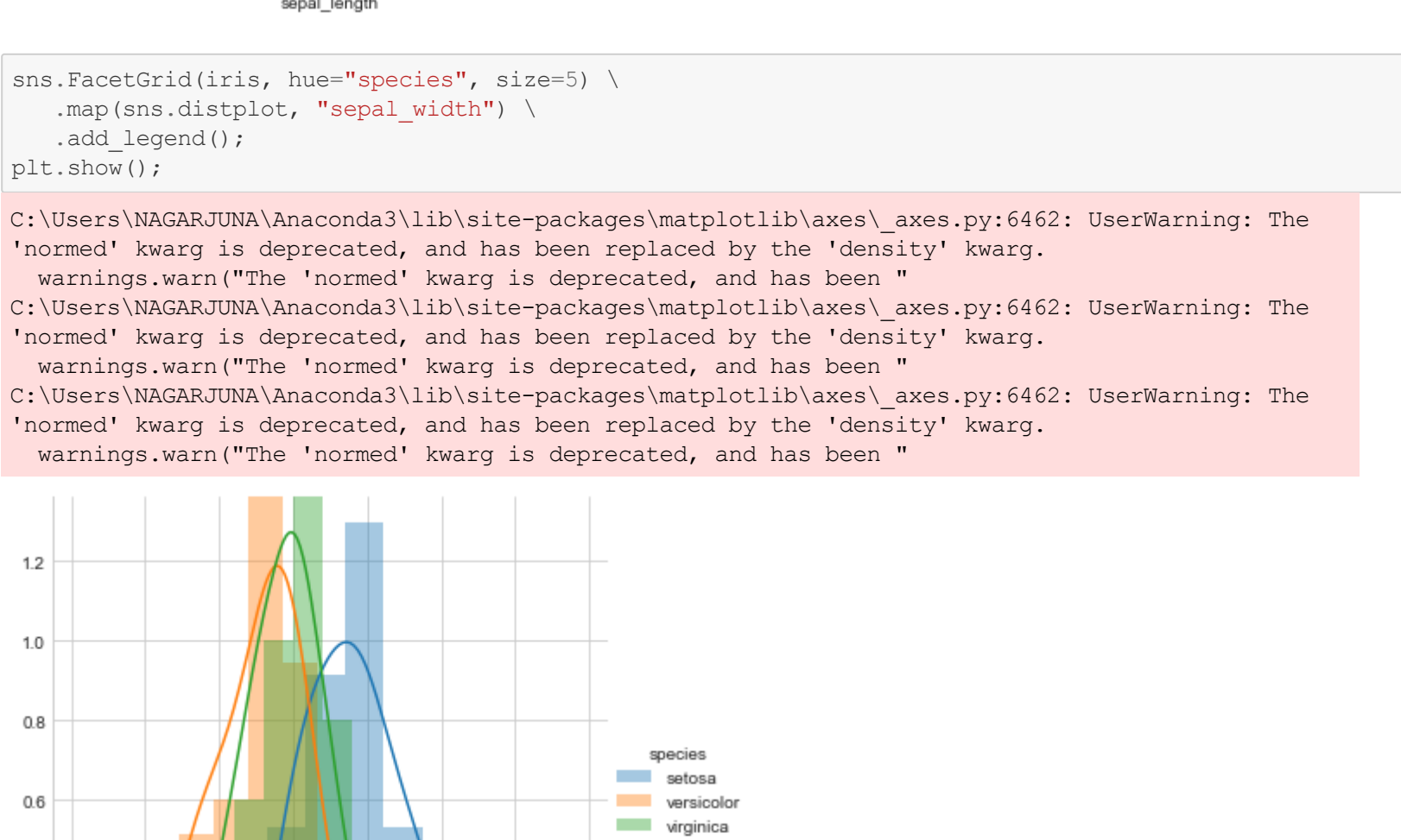
plt.show()

# Disadvantages of 1-D scatter plot: Very hard to make sense as points
# are overlapping a lot.
# Are there better ways of visualizing 1-D scatter plots?
```



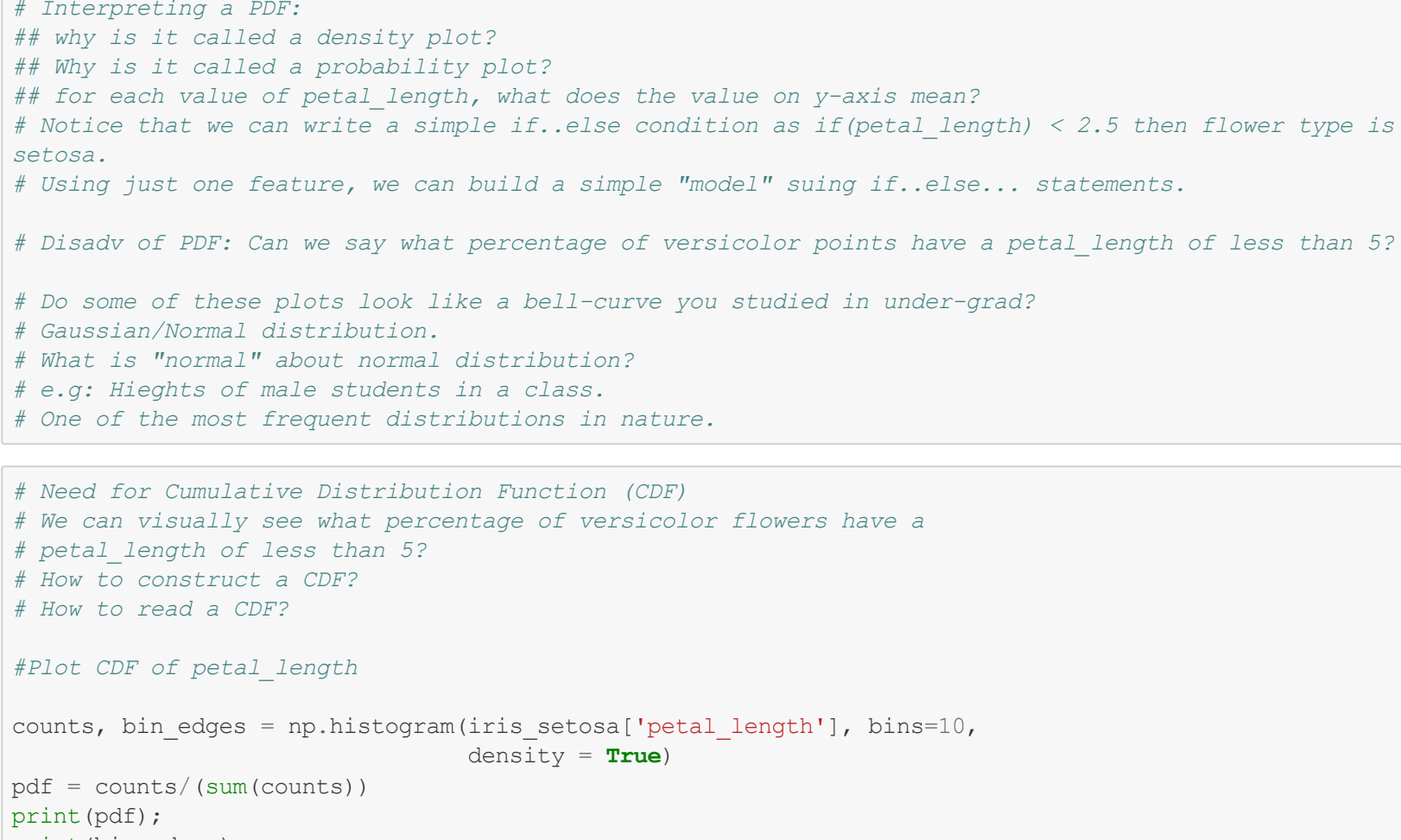
```
In [13]: sns.FacetGrid(iris, hue='species', size=5) \
    .map(sns.distplot, "petal_length") \
    .add_legend() \
    .plt.show()

C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
```



```
In [14]: sns.FacetGrid(iris, hue='species', size=5) \
    .map(sns.distplot, "petal_width") \
    .add_legend() \
    .plt.show()

C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
```



```
In [15]: sns.FacetGrid(iris, hue='species', size=5) \
    .map(sns.distplot, "sepal_length") \
    .add_legend() \
    .plt.show()

C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
```



```
In [16]: sns.FacetGrid(iris, hue='species', size=5) \
    .map(sns.distplot, "sepal_width") \
    .add_legend() \
    .plt.show()

C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
C:\Users\WAGARJUNA\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.",
```



```
In [17]: # Histograms and Probability Density Functions (PDF) using KDE
# How to compute PDFs using counts/frequencies of data points in each window.
# How window width effects the PDF plot.

# Interpreting a PDF:
## Why is it called a density plot?
## Why is it called a probability plot?
## For each value of petal_length, what does the value on y-axis mean?
# Notice that we can write a simple if..else condition as if (petal_mean) < 2.5 then flower type is
setosa.
# Using just one feature, we can build a simple "model" using if..else.. statements.

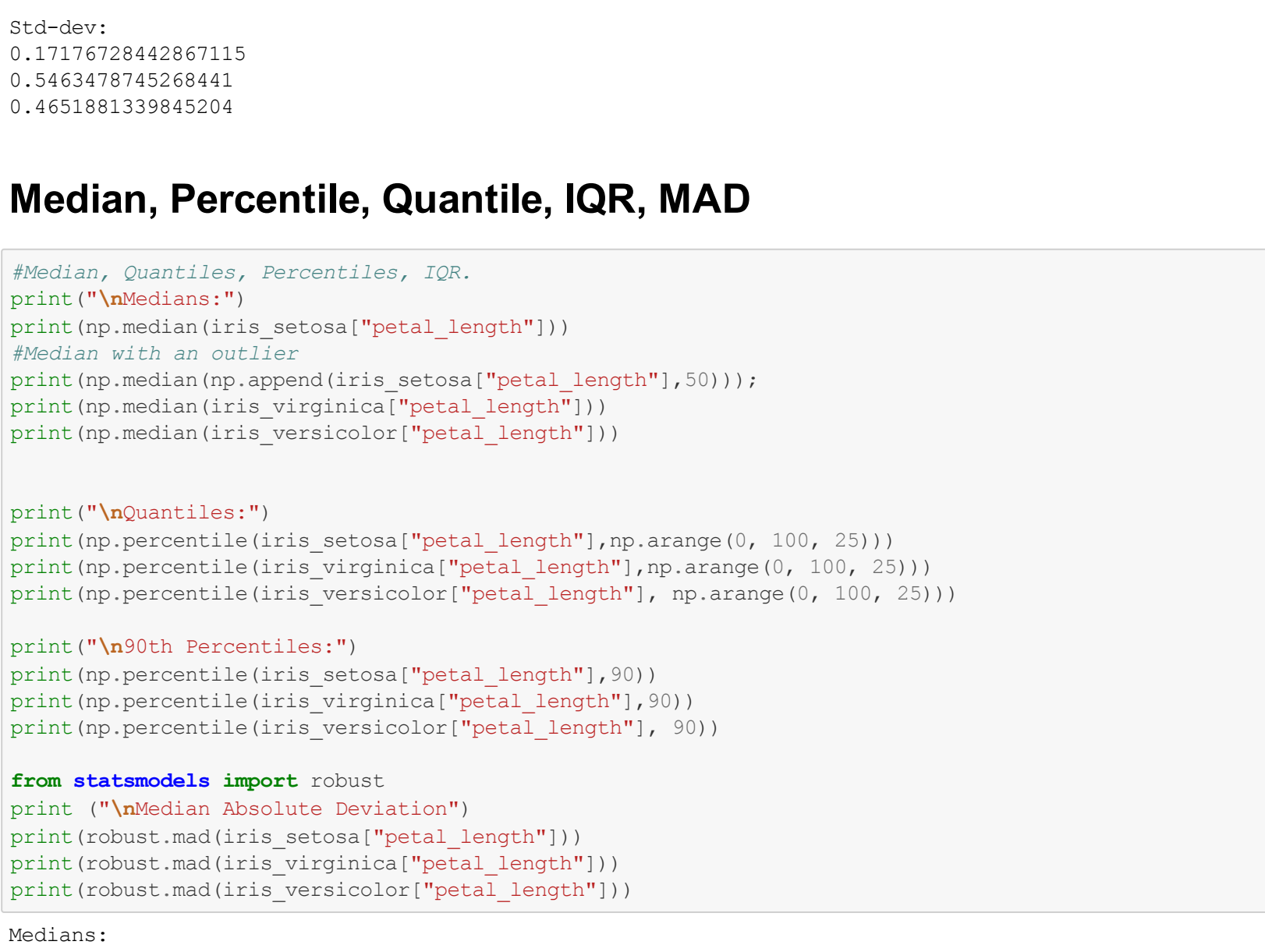
# Disadv of PDF: Can we say what percentage of versicolor points have a petal_length of less than 5?
# Do some of these plots look like a bell-curve you studied in under-grade?
# Gaussian/Normal distribution.
# What is "normal" about normal distribution?
# e.g: Heights of male students in a class.
# One of the most frequent distributions in nature.
```

```
In [18]: # Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolor flowers have a
# petal_length of less than 5?
# How to construct a CDF?
# How to read a CDF?

# Plot CDF of petal_length
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

counts, bin_edges = np.histogram(iris_versicolor['petal_length'], bins=20,
                                density = True)
pdf = counts/(sum(counts))
plt.plot(bin_edges[1:],pdf)

plt.show()
```



```
In [19]: # Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolor flowers have a
# petal_length of less than 1.6?
# How to construct a CDF?
# How to read a CDF?

# Plot CDF of petal_length
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)

# compute CDF
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

plt.show()
```

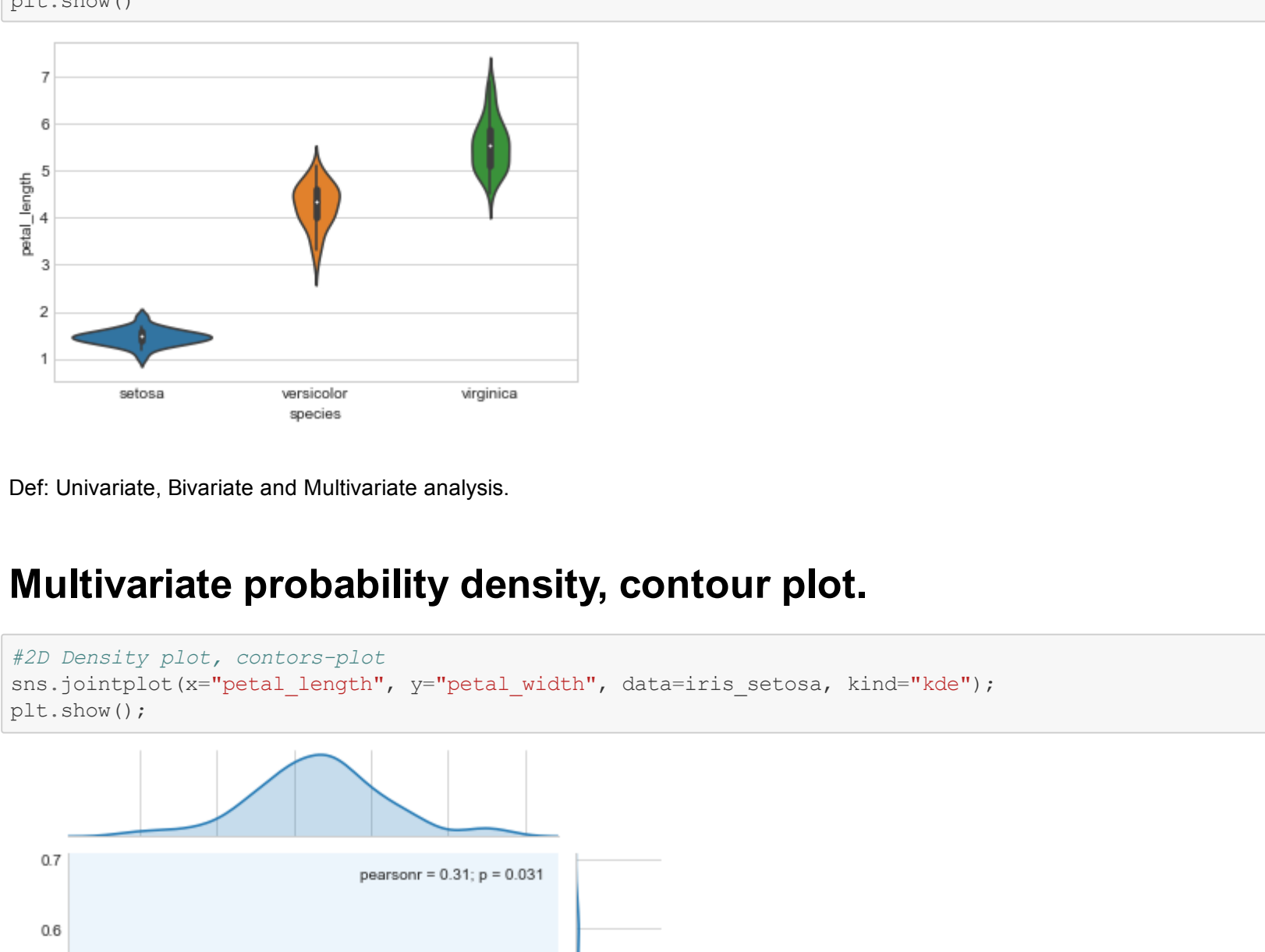


```
In [20]: # Plots of CDF of petal_length for various types of flowers.
# Misclassification error if you use petal_length only.
counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

# virginica
counts, bin_edges = np.histogram(iris_virginica['petal_length'], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

# versicolor
counts, bin_edges = np.histogram(iris_versicolor['petal_length'], bins=10,
                                density = True)
pdf = counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

plt.show()
```



Mean, Variance and Std-dev

```
In [21]: #Mean, Variance, Std-deviation.
print("Means:")
print(np.mean(iris_setosa["petal_length"]))
# Mean with an outlier.
print(np.mean(np.append(iris_setosa["petal_length"], 50)))
print(np.mean(iris_virginica["petal_length"]))
print(np.mean(iris_versicolor["petal_length"]))

print("\nStd-dev:")
print(np.std(iris_setosa["petal_length"]))
print(np.std(iris_virginica["petal_length"]))
print(np.std(iris_versicolor["petal_length"]))

Means:
1.464
2.4156462745098038
5.552
4.262

Std-dev:
0.171767844287115
0.54637874536941
0.6461881339845204
```

Median, Percentile, Quantile, IQR, MAD

```
In [22]: #Median, Quantiles, Percentiles, IQR.
print("\nMedian:")
print(np.median(iris_setosa["petal_length"]))
# Median with an outlier.
print(np.median(np.append(iris_setosa["petal_length"], 50)))
print(np.median(iris_virginica["petal_length"]))
print(np.median(iris_versicolor["petal_length"]))

print("\nQuantiles:")
print(np.percentile(iris_setosa["petal_length"], np.arange(0, 100, 25)))
print(np.percentile(iris_virginica["petal_length"], np.arange(0, 100, 25)))
print(np.percentile(iris_versicolor["petal_length"], np.arange(0, 100, 25)))

print("\n30th Percentiles:")
print(np.percentile(iris_setosa["petal_length"], 30))
print(np.percentile(iris_virginica["petal_length"], 30))
print(np.percentile(iris_versicolor["petal_length"], 30))

from statsmodels import robust
print("\nMedian Absolute Deviation")
print(robust.mad(iris_setosa["petal_length"]))
print(robust.mad(iris_virginica["petal_length"]))
print(robust.mad(iris_versicolor["petal_length"]))

Medians:
1.5
1.5
5.55
4.35

Quantiles:
[1. 4. 5. 1.5 1.575]
[4.5 5.1 5.55 5.875]
[3. 4. 4.35 4.6 ]

90th Percentiles:
1.7
6.3100000000000005
4.8

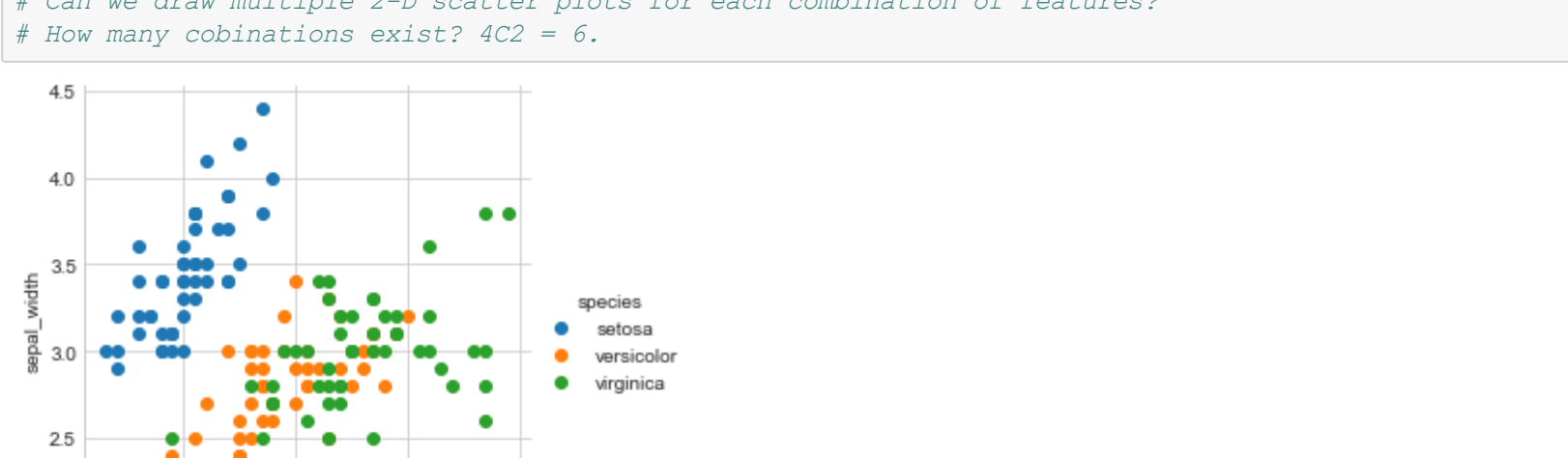
Median Absolute Deviation
0.1482602218505631
0.4671709883275113
0.518197764769602
```

Box plot and Whiskers

```
In [23]: #Box-plot with whiskers: another method of visualizing the 1-D scatter plot more intuitively.
# The Concept of median, percentile, quantile.
# How to draw the box in a box-plot?
# How to draw whiskers: (no standard way) Could use min and max or use other complex statistical: two
hinges
# IQR-like idea.

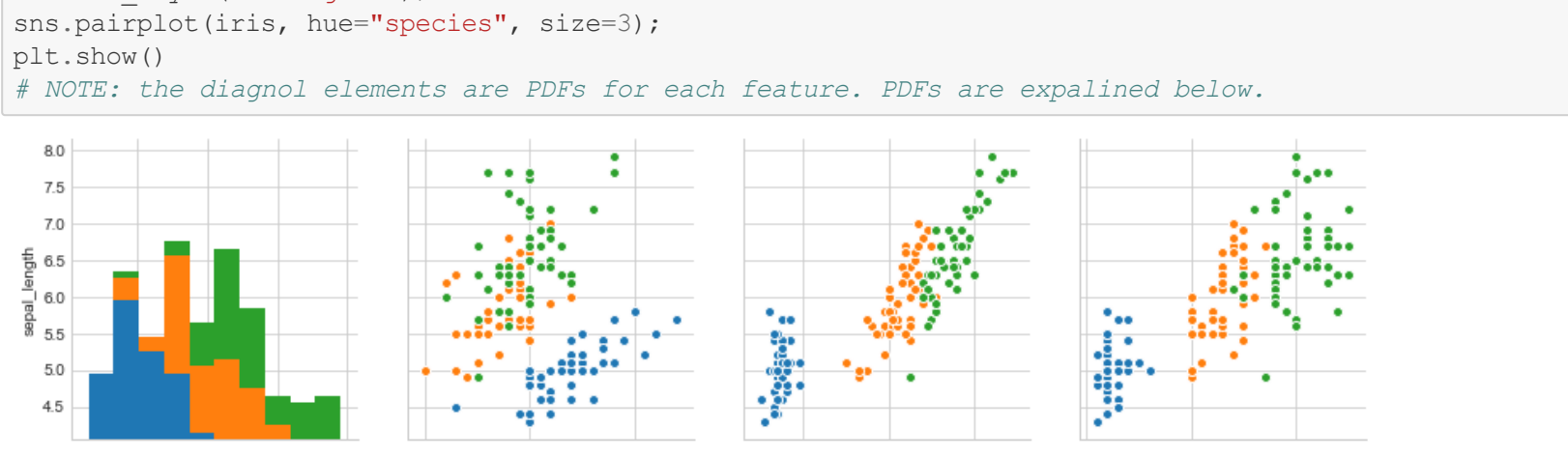
# NOTE: IN the plot below, a technique call Inter-quartile range is used in plotting the whiskers.
# Whiskers in the plot below donot correspond to the min and max values.

# Box-plot can be visualized as a PDF on the side-ways.
sns.boxplot(x='species', y='petal_length', data=iris)
plt.show()
```



Violin plots

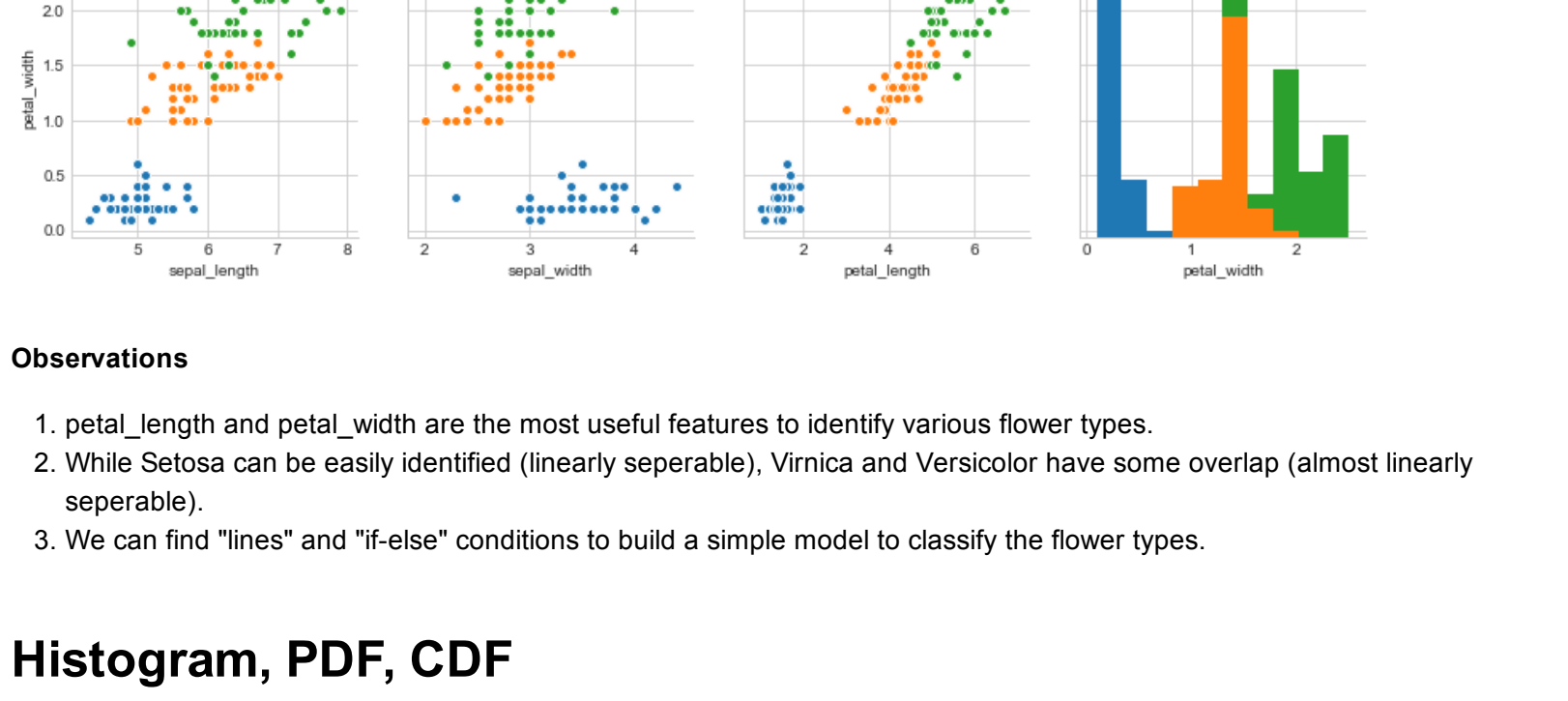
```
In [24]: # A violin plot combines the benefits of the previous two plots
# and simplifies them
# Denser regions of the data are fatter, and sparser ones thinner
# In a violin plot
sns.violinplot(x='species', y='petal_length', data=iris, size=8)
plt.show()
```



Def: Univariate, Bivariate and Multivariate analysis.

Multivariate probability density, contour plot.

```
In [26]: #2D Density plot, contours-plot
sns.jointplot(x="petal_length", y="petal_width", data=iris_setosa, kind="kde")
plt.show()
```



- Observations:**
- Using Petals over Sepal for training the data gives a much better accuracy. This was expected as we saw in the heatmap above that the correlation between the Sepal Width and Length was very low whereas the correlation between Petal Width and Length was very high.