

## Iris Flower dataset

```
In [2]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

'''downlaod iris.csv from https://raw.githubusercontent.com/uiuc-cse/data-fa14/gh-pages/data/iris.csv'''
#Load Iris.csv into a pandas DataFrame.
iris = pd.read_csv("iris.csv")

In [3]: # (Q) how many data-points and features?
print (iris.shape)

(150, 5)

In [4]: #(Q) What are the column names in our dataset?
print (iris.columns)

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

In [5]: #(Q) How many data points for each class are present?
#(or) How many flowers for each species are present?

iris["species"].value_counts()
# balanced-dataset vs imbalanced datasets
#Iris is a balanced dataset as the number of data points for every class is 50.

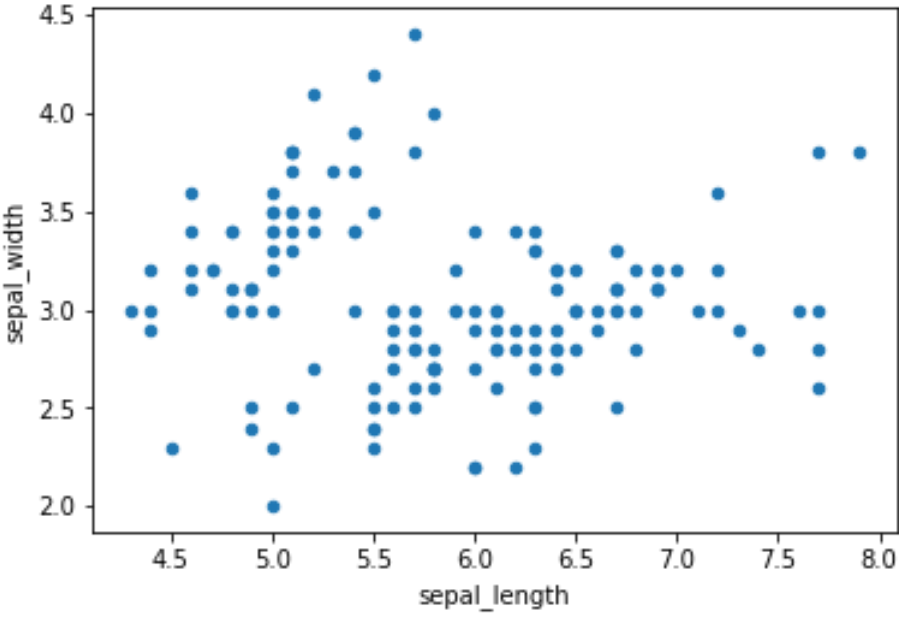
Out[5]: virginica      50
setosa      50
versicolor  50
Name: species, dtype: int64
```

## 2-D Scatter Plot

```
In [6]: #2-D scatter plot:
#ALWAYS understand the axis: labels and scale.

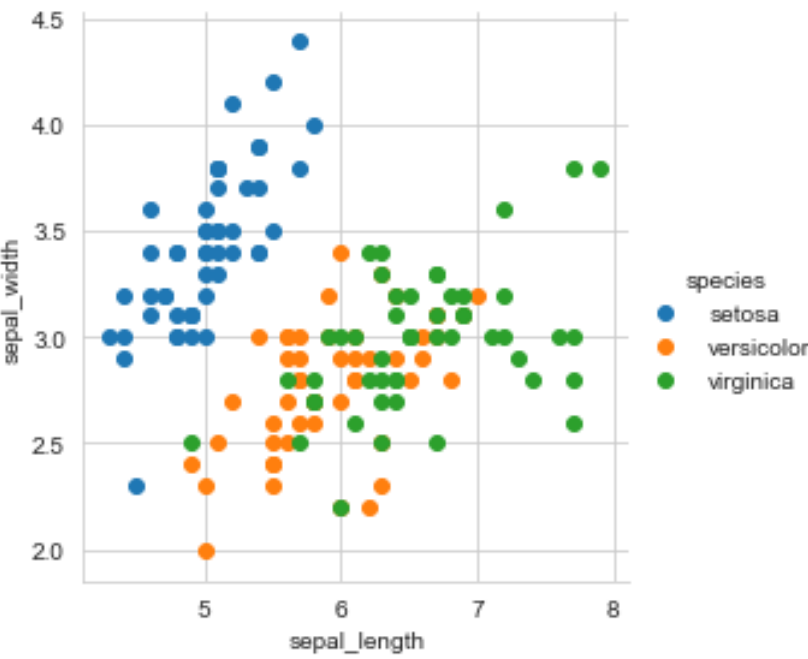
iris.plot(kind='scatter', x='sepal_length', y='sepal_width') ;
plt.show()

#cannot make much sense out it.
#What if we color the points by thier class-label/flower-type.
```



```
In [8]: # 2-D Scatter plot with color-coding for each flower type/class.
# Here 'sns' corresponds to seaborn.
sns.set_style("whitegrid");
sns.FacetGrid(iris, hue="species", height=4) \
    .map(plt.scatter, "sepal_length", "sepal_width") \
    .add_legend();
plt.show();

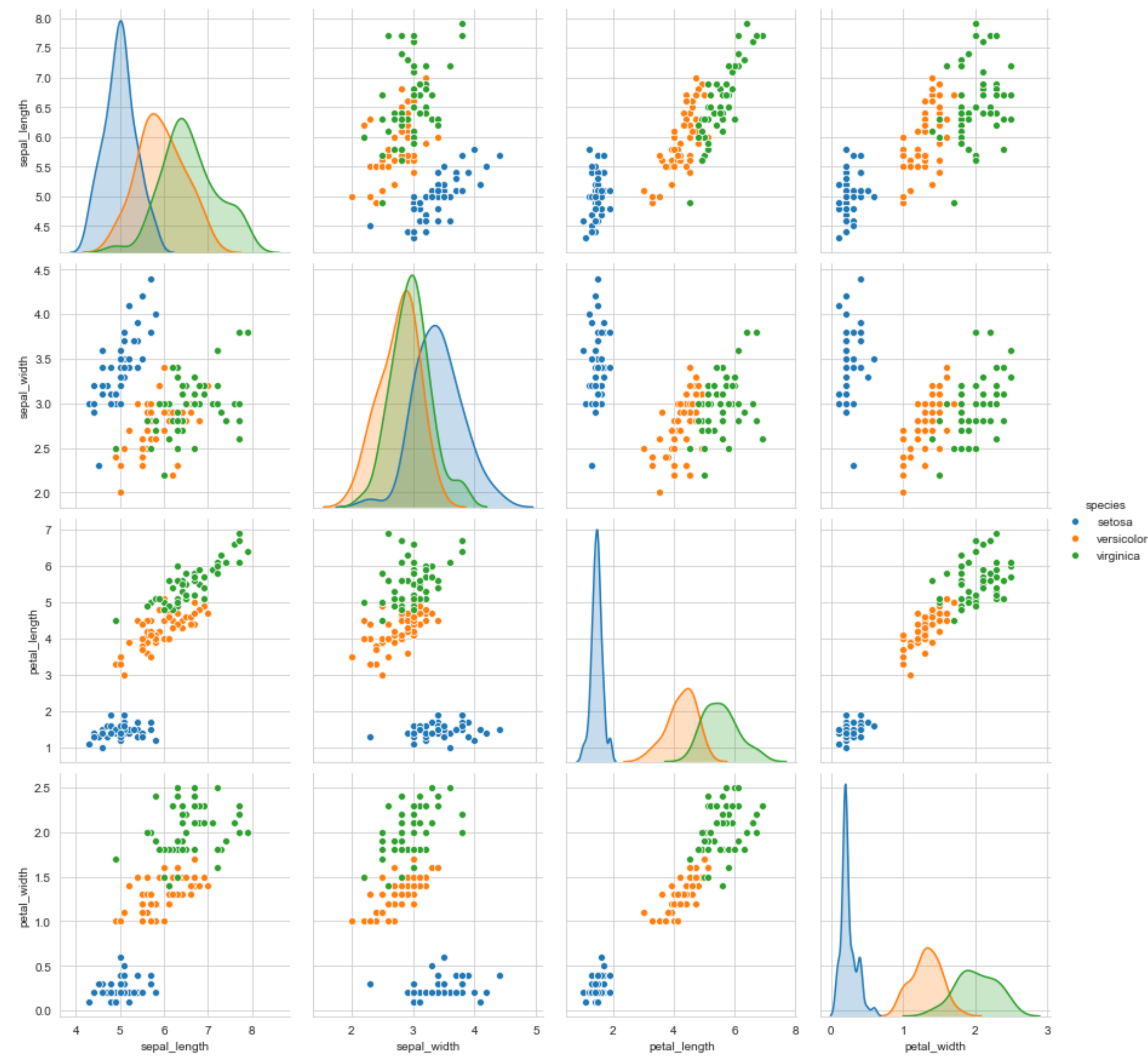
# Notice that the blue points can be easily seperated
# from red and green by drawing a line.
# But red and green data points cannot be easily seperated.
# Can we draw multiple 2-D scatter plots for each combination of features?
# How many cobinations exist? 4C2 = 6.
```



- Observation(s):**
- 1. Using sepal\_length and sepal\_width features, we can distinguish Setosa flowers from others.
  - 2. Seperating Versicolor from Viginica is much harder as they have considerable overlap.

**Pair-plot**

```
In [9]: # 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", height=3);
plt.show()
# NOTE: the diagnol elements are PDFs for each feature. PDFs are expalined below.
```



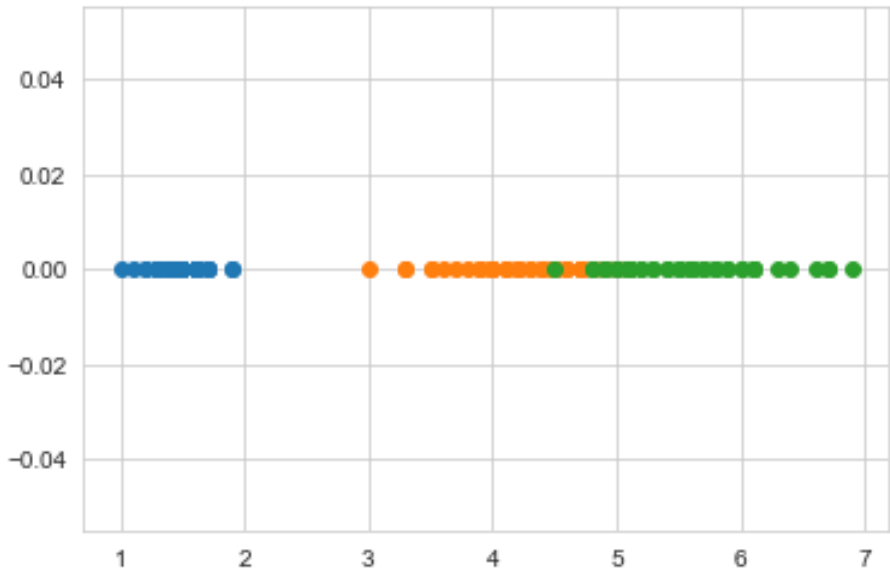
**Observations**

- 1. petal\_length and petal\_width are the most useful features to identify various flower types.
- 2. While Setosa can be easily identified (linearly seperable), Virnica and Versicolor have some overlap (almost linearly seperable).
- 3. We can find "lines" and "if-else" conditions to build a simple model to classify the flower types.

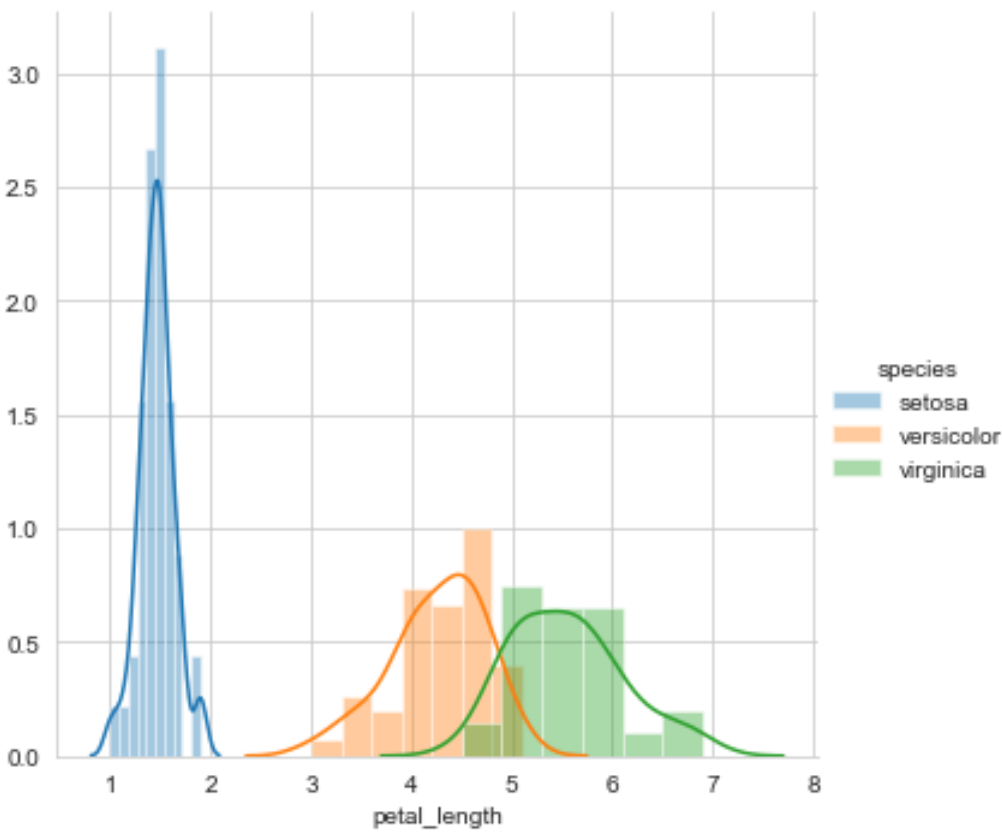
## Histogram, PDF, CDF

```
In [10]: # 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"])
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')

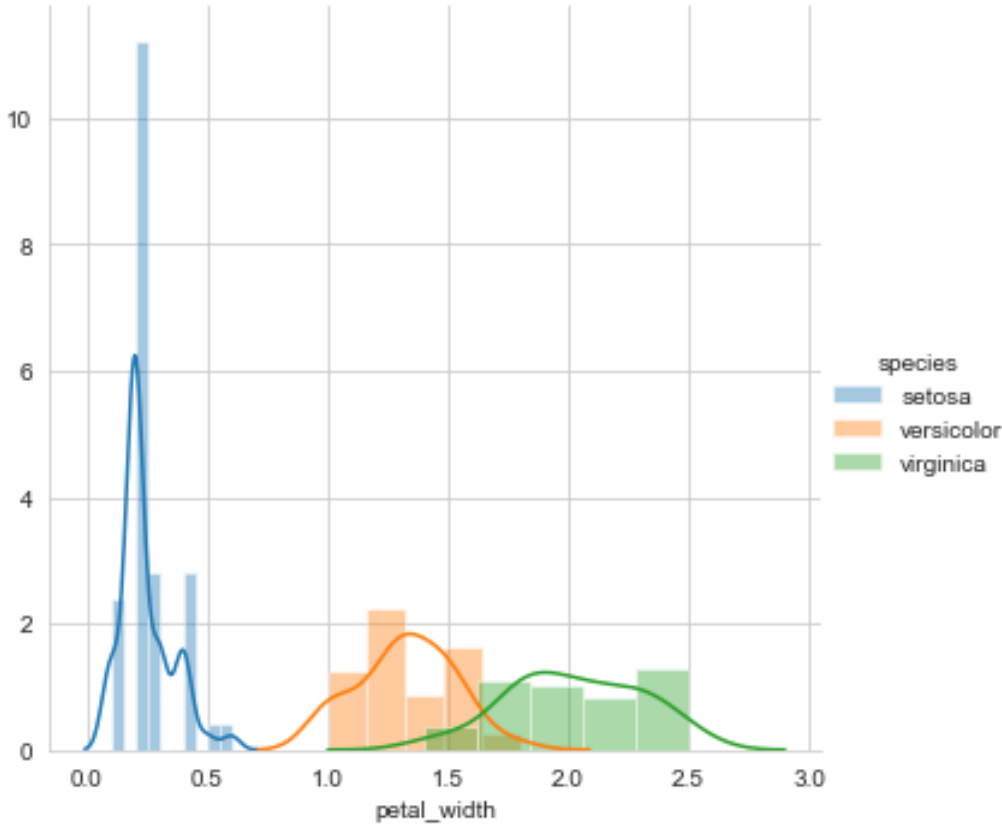
plt.show()
```



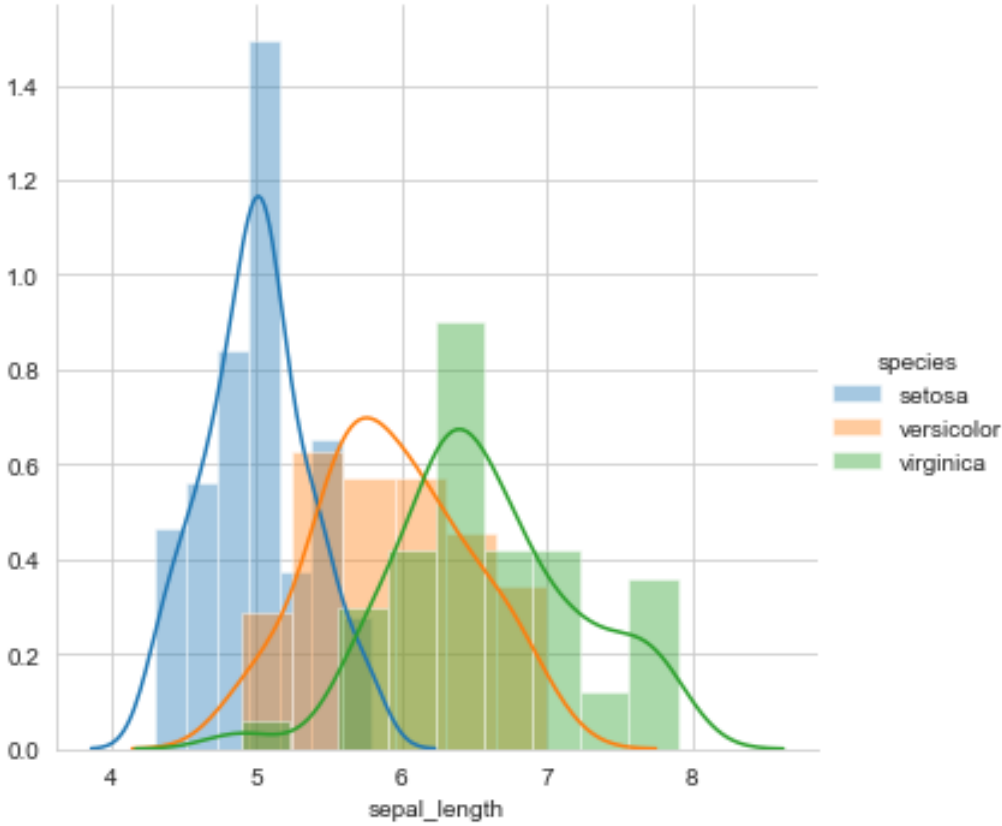
```
In [11]: sns.FacetGrid(iris, hue="species", height=5) \
    .map(sns.distplot, "petal_length") \
    .add_legend();
plt.show();
```



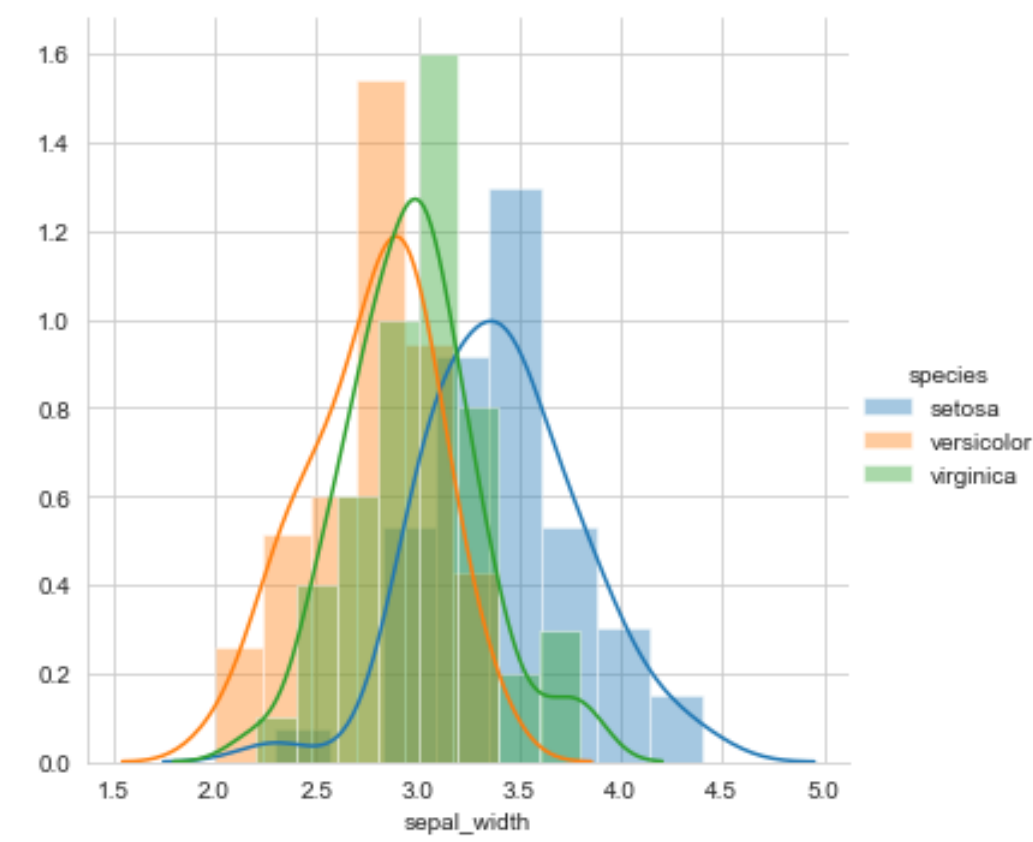
```
In [12]: sns.FacetGrid(iris, hue="species", height=5) \
        .map(sns.distplot, "petal_width") \
        .add_legend();
plt.show();
```



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



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



```
In [15]: # How to compute PDFs using counts/frequencies of data points in each window.
# How window width effects the PDF plot.

# Interpreting a PDF:

## 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_length) < 2.5 then flower type is setosa.
# Using just one feature, we can build a simple "model" suing if..else... statements.

# Can we say what percentage of versicolor points have a petal_length of less than 5?
```

In [16]: *# We can visually see what percentage of versicolor flowers have a*  
*# petal\_length of less than 5?*

```
#Plot CDF of petal_length

counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=10,
                                density = True)

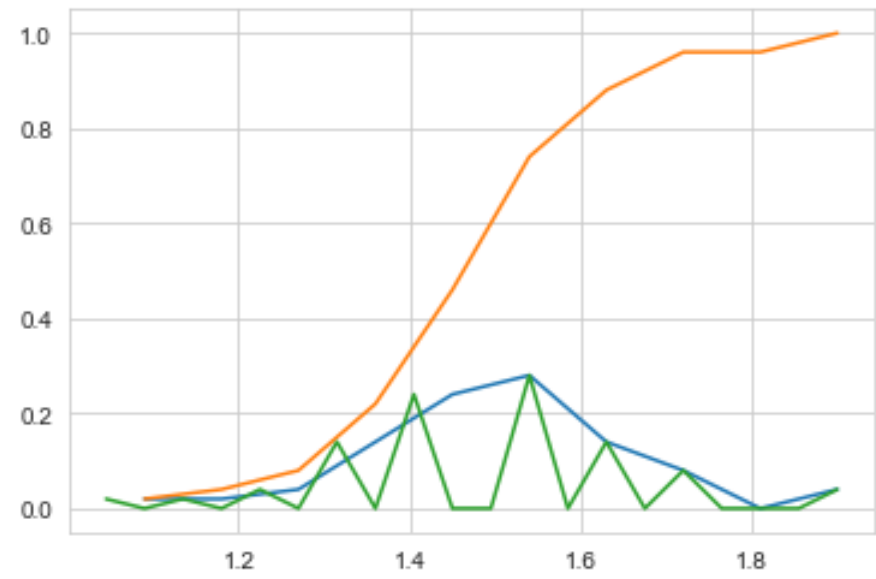
pdf = counts/(sum(counts))
print(pdf);
print(bin_edges);
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf);
plt.plot(bin_edges[1:], cdf)

counts, bin_edges = np.histogram(iris_setosa['petal_length'], bins=20,
                                density = True)

pdf = counts/(sum(counts))
plt.plot(bin_edges[1:],pdf);

plt.show();
```

[0.02 0.02 0.04 0.14 0.24 0.28 0.14 0.08 0. 0.04]  
[1. 1.09 1.18 1.27 1.36 1.45 1.54 1.63 1.72 1.81 1.9 ]



```
In [17]: # We can visually see what percentage of versicolor flowers have a
# petal_length of less than 1.6?
```

```
#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
```

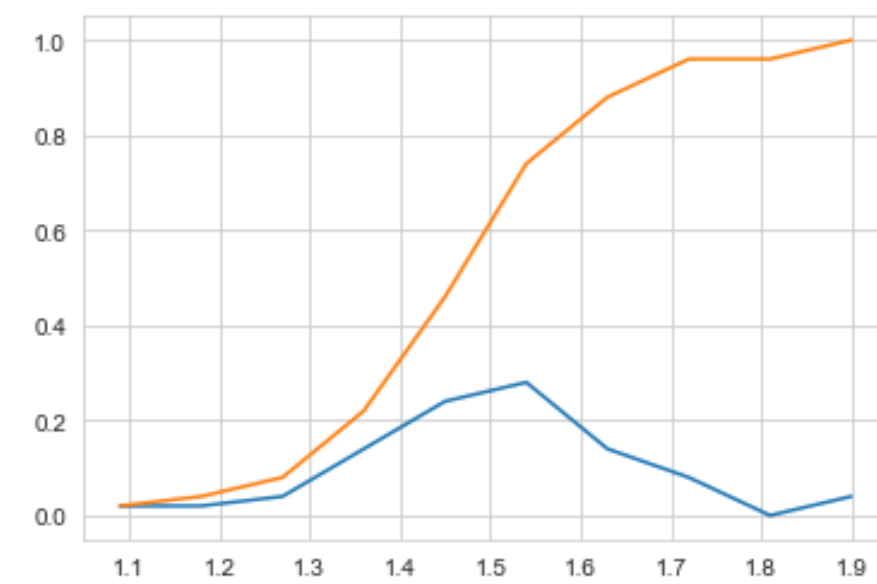
```
cdf = np.cumsum(pdf)
```

```
plt.plot(bin_edges[1:],pdf)
```

```
plt.plot(bin_edges[1:], cdf)
```

```
plt.show();
```

```
[0.02 0.02 0.04 0.14 0.24 0.28 0.14 0.08 0.    0.04]
[1.    1.09 1.18 1.27 1.36 1.45 1.54 1.63 1.72 1.81 1.9 ]
```





```
In [18]: # 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)

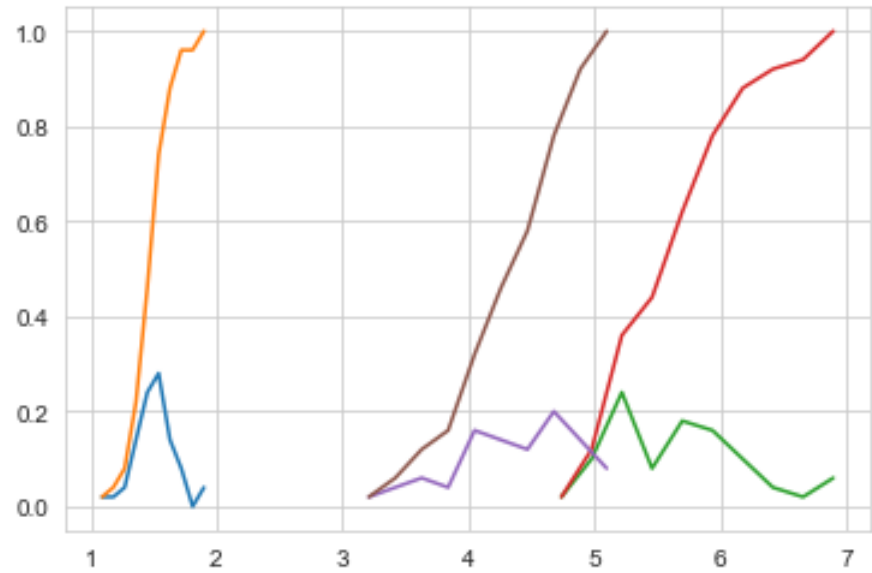
pdf = counts/(sum(counts))
print(pdf);
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:], cdf)

#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();
```

```
[0.02 0.02 0.04 0.14 0.24 0.28 0.14 0.08 0.  0.04]
[1.   1.09 1.18 1.27 1.36 1.45 1.54 1.63 1.72 1.81 1.9 ]
[0.02 0.1  0.24 0.08 0.18 0.16 0.1  0.04 0.02 0.06]
[4.5  4.74 4.98 5.22 5.46 5.7  5.94 6.18 6.42 6.66 6.9 ]
[0.02 0.04 0.06 0.04 0.16 0.14 0.12 0.2  0.14 0.08]
[3.   3.21 3.42 3.63 3.84 4.05 4.26 4.47 4.68 4.89 5.1 ]
```



Mean, Variance and Std-dev

```
In [19]: #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
2.4156862745098038
5.552
4.26

Std-dev:
0.17176728442867115
0.5463478745268441
0.4651881339845204
```

Median, Percentile, Quantile, IQR, MAD

```
In [20]: #Median, Quantiles, Percentiles, IQR.
print("\nMedians:")
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("\n90th Percentiles:")
print(np.percentile(iris_setosa["petal_length"],90))
print(np.percentile(iris_virginica["petal_length"],90))
print(np.percentile(iris_versicolor["petal_length"], 90))

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

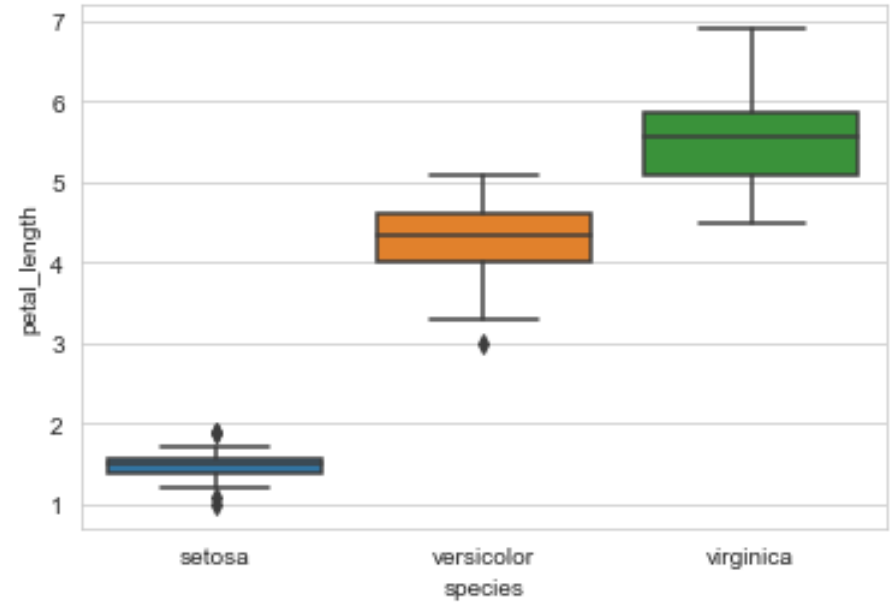
Quantiles:  
[1.    1.4    1.5    1.575]  
[4.5   5.1    5.55  5.875]  
[3.    4.    4.35 4.6 ]

90th Percentiles:  
1.7  
6.3100000000000005  
4.8

Median Absolute Deviation  
0.14826022185056031  
0.6671709983275211  
0.5189107764769602

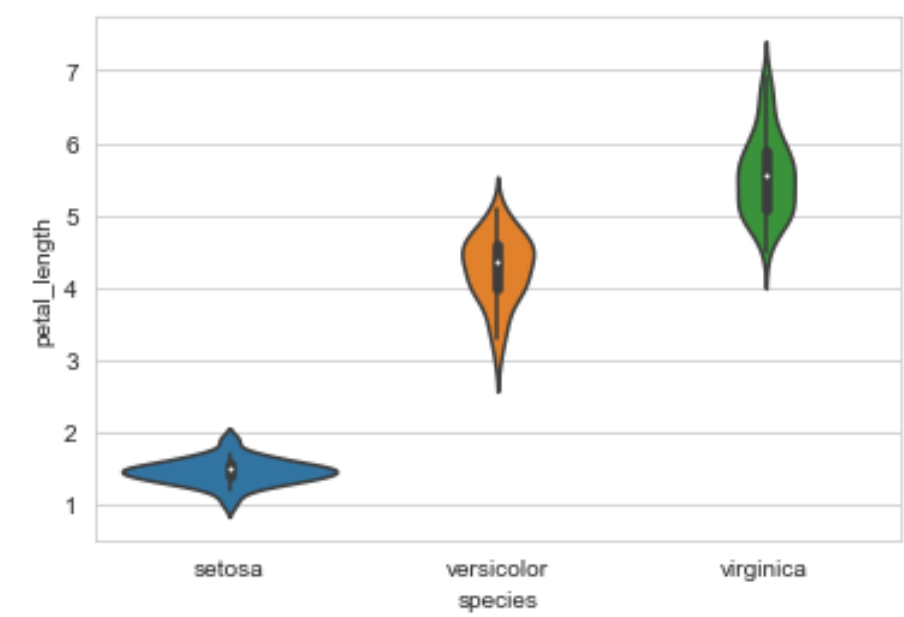
Box plot and Whiskers

```
In [21]: sns.boxplot(x='species',y='petal_length', data=iris)
plt.show()
```



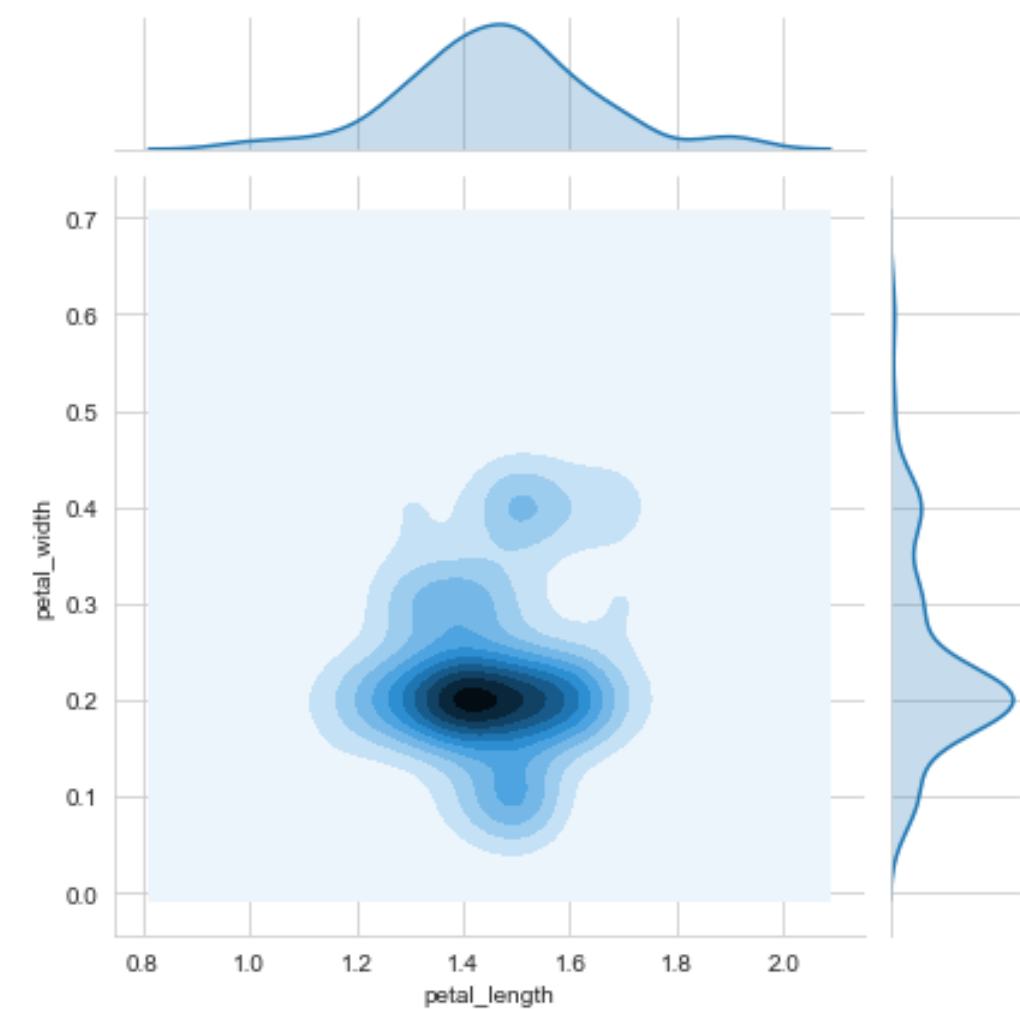
## Violin plots

```
In [22]: sns.violinplot(x="species", y="petal_length", data=iris, size=8)
plt.show()
```



## Multivariate probability density, contour plot.

```
In [23]: #2D Density plot, contours-plot
sns.jointplot(x="petal_length", y="petal_width", data=iris_setosa, kind="kde");
plt.show();
```



```
In [25]: iris_virginica_SW = iris_virginica.iloc[:,1]
iris_versicolor_SW = iris_versicolor.iloc[:,1]
```

```
In [26]: from scipy import stats
stats.ks_2samp(iris_virginica_SW, iris_versicolor_SW)
```

Out[26]: Ks\_2sampResult(statistic=0.26, pvalue=0.06779471096995852)

```
In [27]: x = stats.norm.rvs(loc=0.2, size=10)
stats.kstest(x, 'norm')
```

Out[27]: KstestResult(statistic=0.24891960150941622, pvalue=0.495575609780464)

```
In [28]: x = stats.norm.rvs(loc=0.2, size=100)
stats.kstest(x, 'norm')
```

```
Out[28]: KstestResult(statistic=0.14332186743894548, pvalue=0.029573824992372134)
```

```
In [29]: x = stats.norm.rvs(loc=0.2, size=1000)
stats.kstest(x, 'norm')
```

```
Out[29]: KstestResult(statistic=0.1025760241341469, pvalue=1.2960861619016006e-09)
```