```
In [43]: import pandas as pd
         import numpy as np
         import seaborn as sns
         import matplotlib.pyplot as plt
         import warnings
```

```
In [44]: from sklearn import metrics
         from sklearn.neighbors import KNeighborsClassifier
         from sklearn.linear_model import LogisticRegression
         from sklearn.model_selection import train_test_split
```

# Loading the data

```
In [45]: iris = pd.read_csv("iris.csv")
```

```
In [46]: print(iris.shape)
```

```
(150, 5)
```

```
In [47]: iris.head()
```

Out[47]:

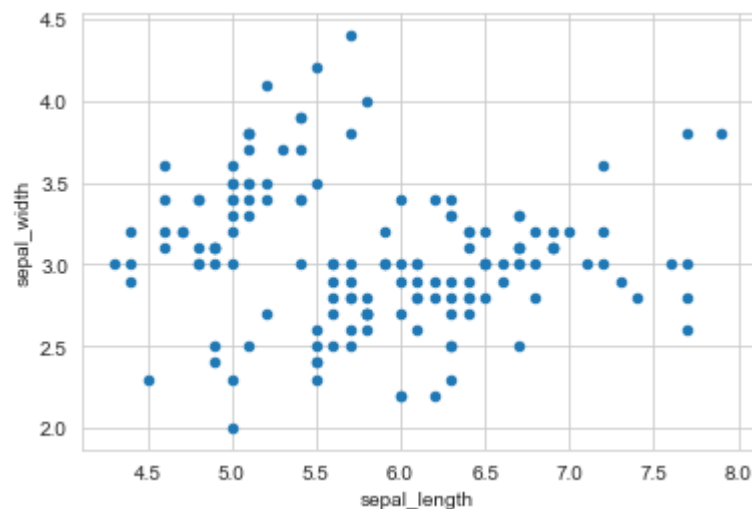|   | sepal_length | sepal_width | petal_length | petal_width | species |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

```
In [48]: iris["species"].value_counts()
```

```
Out[48]: versicolor    50
         setosa        50
         virginica     50
         Name: species, dtype: int64
```

This clearly indicates that the dataset iris is clealy a balanced dataset. And each of the species contains the same number of datapoints
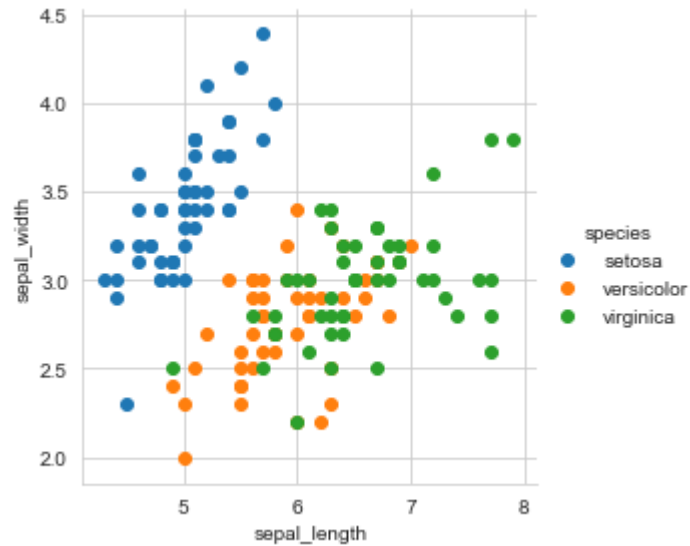
## 2DD Scatter plot

```
In [49]: iris.plot(kind='scatter', x='sepal_length', y='sepal_width') ;
         plt.show()
```



this does not make sence out of it let us visualize more i.e use colour for each class so that we can define the things better

```
In [50]: 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.
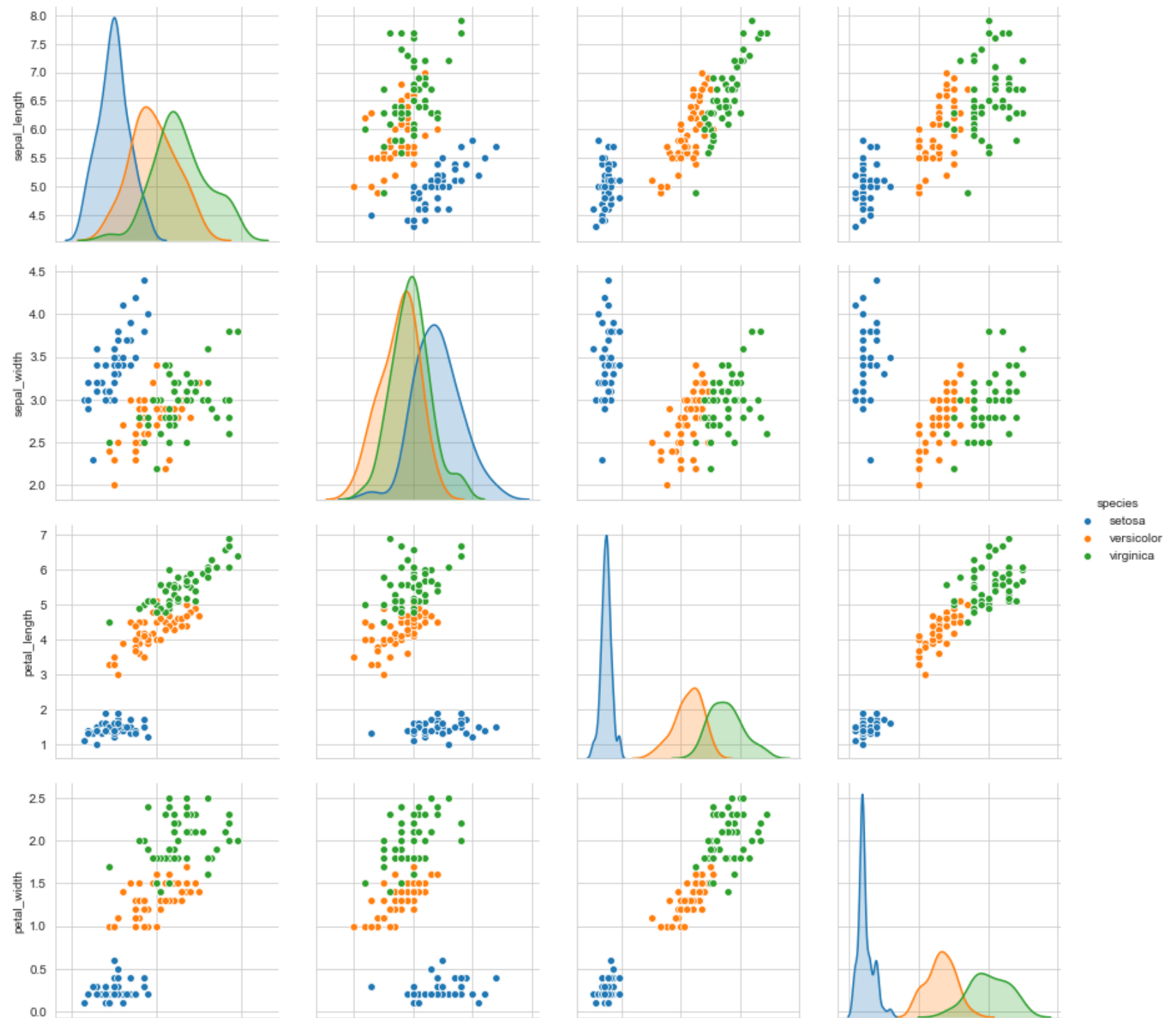
## Observation
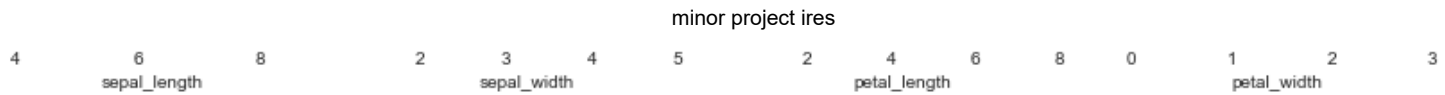
1: Using sepal_length and sepal_width features, we can distinguish Setosa flowers from others.

2: Seperating Versicolor from Viginica is much harder as they have considerable overlap.

# 3D Scatter plot

# Pair plot

In [51]:
```python
plt.close();
sns.set_style("whitegrid");
sns.pairplot(iris, hue="species", height=3);
plt.show()
```

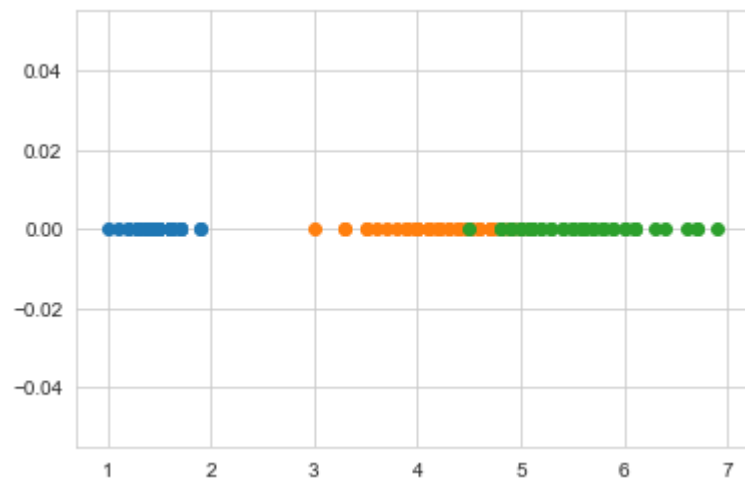|          |          |          |          |
|----------|----------|----------|----------|
| 4        6        8 | 2     3    4     5 | 2       4    6    8 | 0       1      2      3 |
| sepal_length | sepal_width | petal_length | petal_width |

**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 [52]: 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()
```
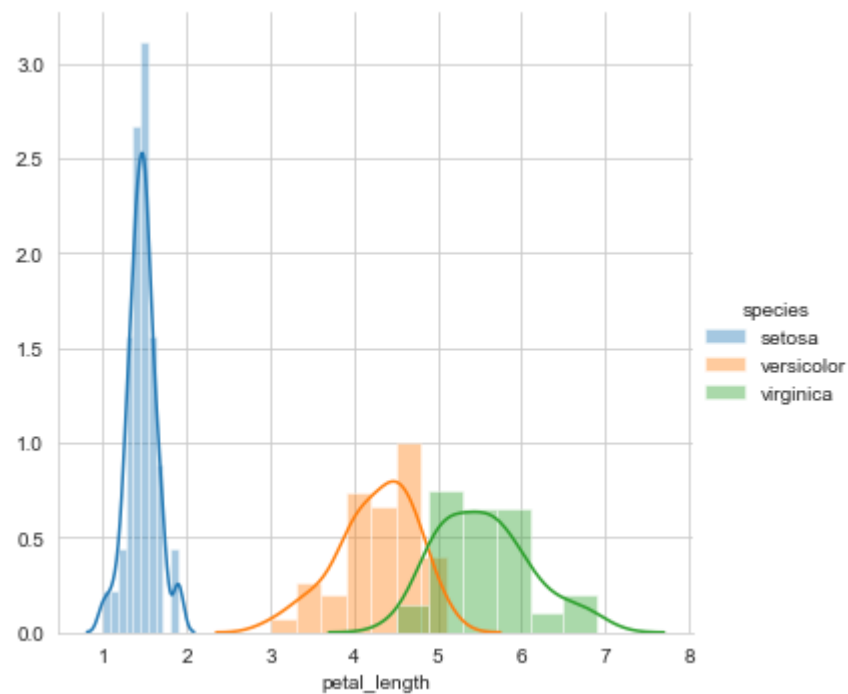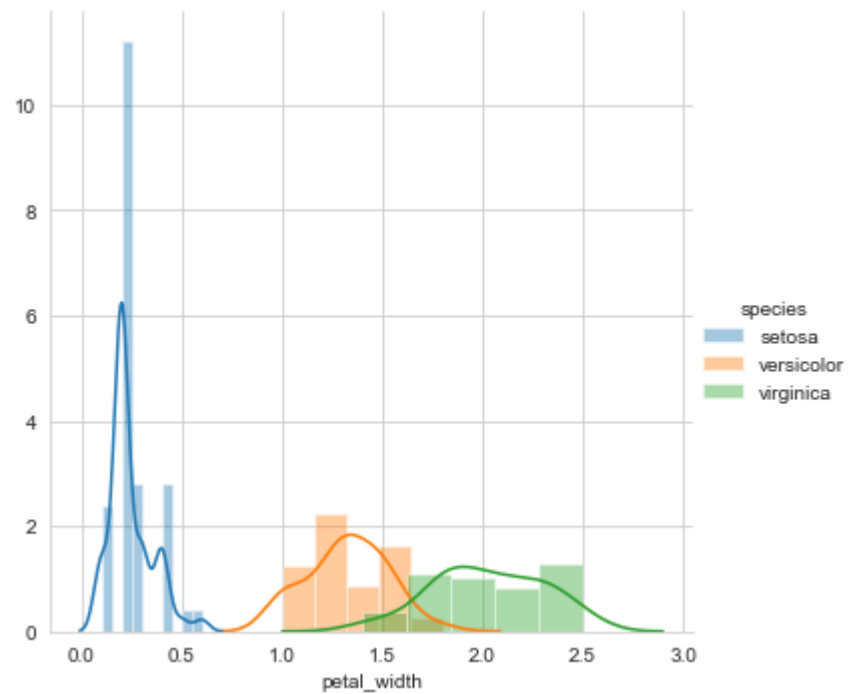
Here the points are overlapping a lot.

```
In [53]: sns.FacetGrid(iris, hue="species", size=5) \
            .map(sns.distplot, "petal_length") \
            .add_legend();
         plt.show();
```
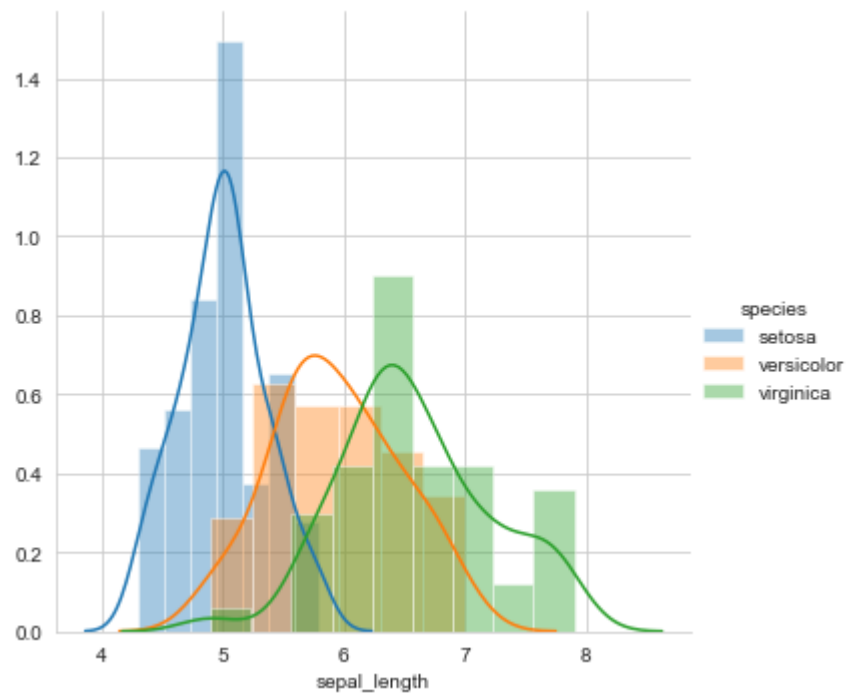
C:\Users\SUNNY\Anaconda3\lib\site-packages\seaborn\axisgrid.py:230: UserWarning: The `size` paramter has been renamed to `height`; please update your code.
  warnings.warn(msg, UserWarning)

In [54]:
```python
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "petal_width") \
    .add_legend();
plt.show();
```

In [55]:
```python
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "sepal_length") \
    .add_legend();
plt.show();
```
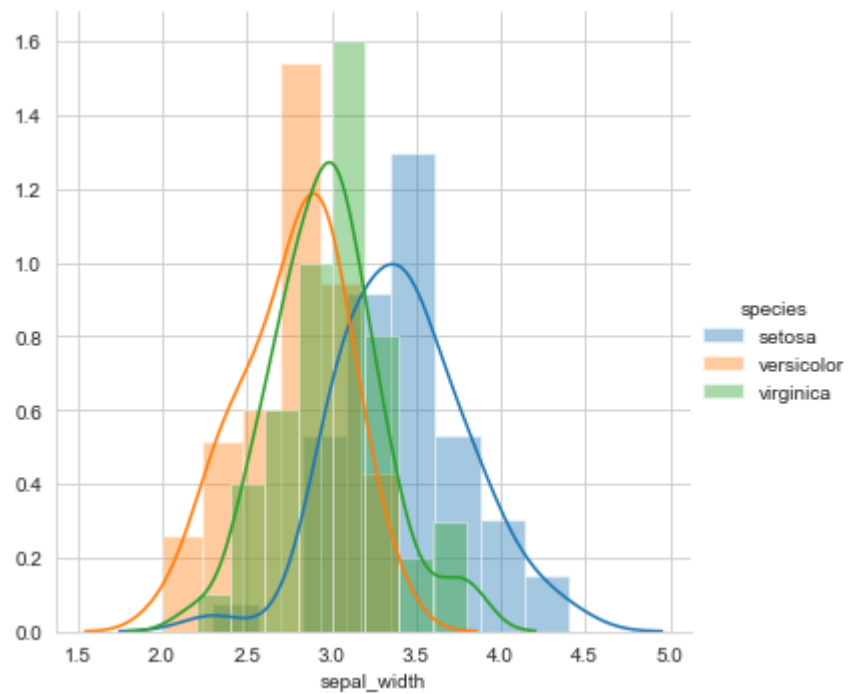
In [56]:
```python
sns.FacetGrid(iris, hue="species", size=5) \
    .map(sns.distplot, "sepal_width") \
    .add_legend();
plt.show();
```

In [57]:
```python
# Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolor flowers have a

#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 [58]:
```python
# Need for Cumulative Distribution Function (CDF)
# We can visually see what percentage of versicolor flowers have a

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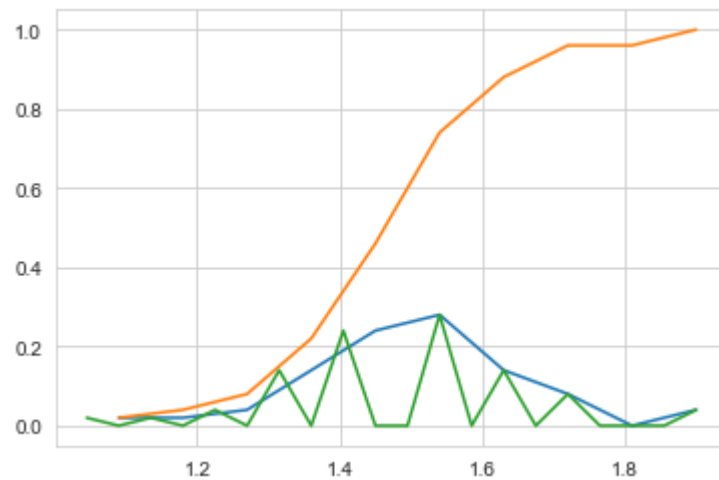
```
[0.02 0.02 0.04 0.14 0.24 0.28 0.14 0.08 0.   0.04]
[1.   1.09 1.18 1.27 1.36 1.45 1.54 1.63 1.72 1.81 1.9 ]
```

In [59]:

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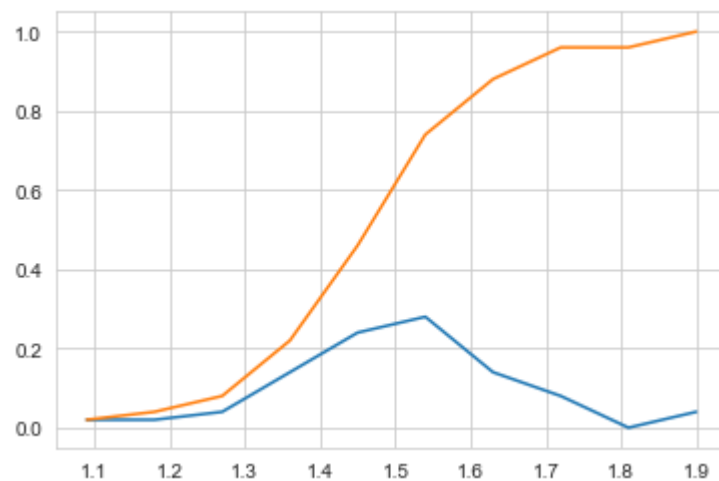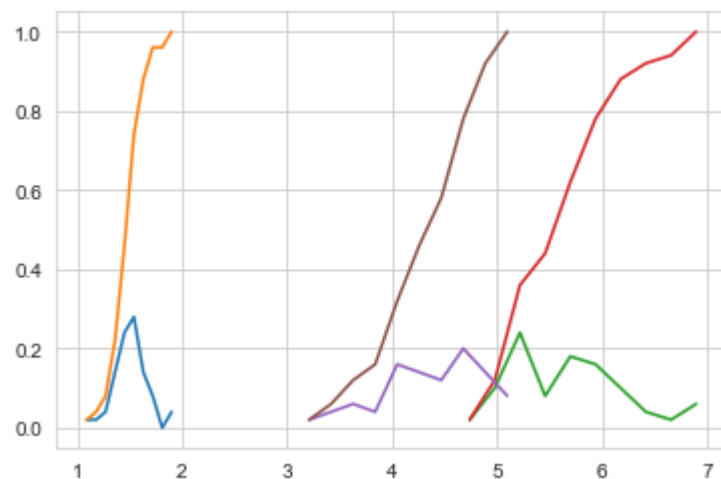
```
[0.02 0.02 0.04 0.14 0.24 0.28 0.14 0.08 0.   0.04]
[1.   1.09 1.18 1.27 1.36 1.45 1.54 1.63 1.72 1.81 1.9 ]
[0.02 0.1  0.24 0.08 0.18 0.16 0.1  0.04 0.02 0.06]
[4.5  4.74 4.98 5.22 5.46 5.7  5.94 6.18 6.42 6.66 6.9 ]
[0.02 0.04 0.06 0.04 0.16 0.14 0.12 0.2  0.14 0.08]
[3.   3.21 3.42 3.63 3.84 4.05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



# Mean, Variance and Std-dev

```
In [60]: #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 [61]:  #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
```
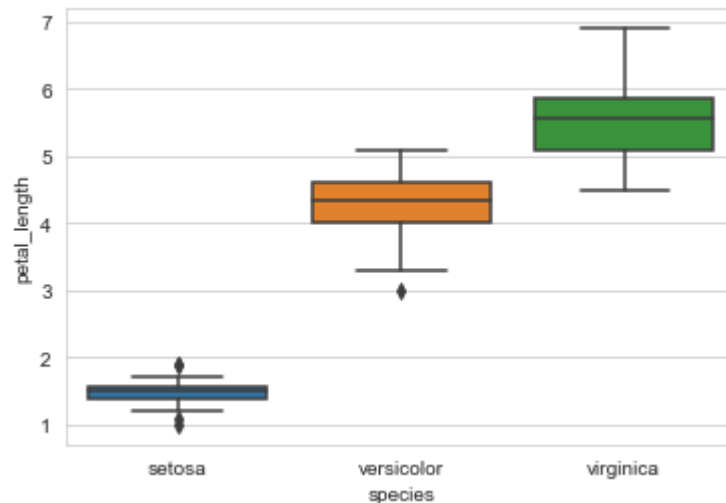
# Box plot and Whiskers

In [62]: 
```
#Box-plot with whiskers: another method of visualizing the  1-D scatter plot more intuitivey.
# The Concept of median, percentile, quantile.


# IN the plot below, a technique call inter-quartile range is used in plotting the whiskers.
#Whiskers in the plot below donot correposnd 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()
```
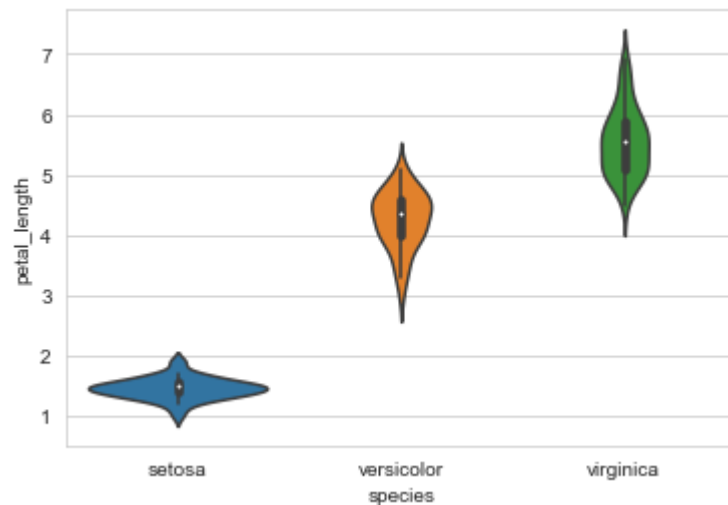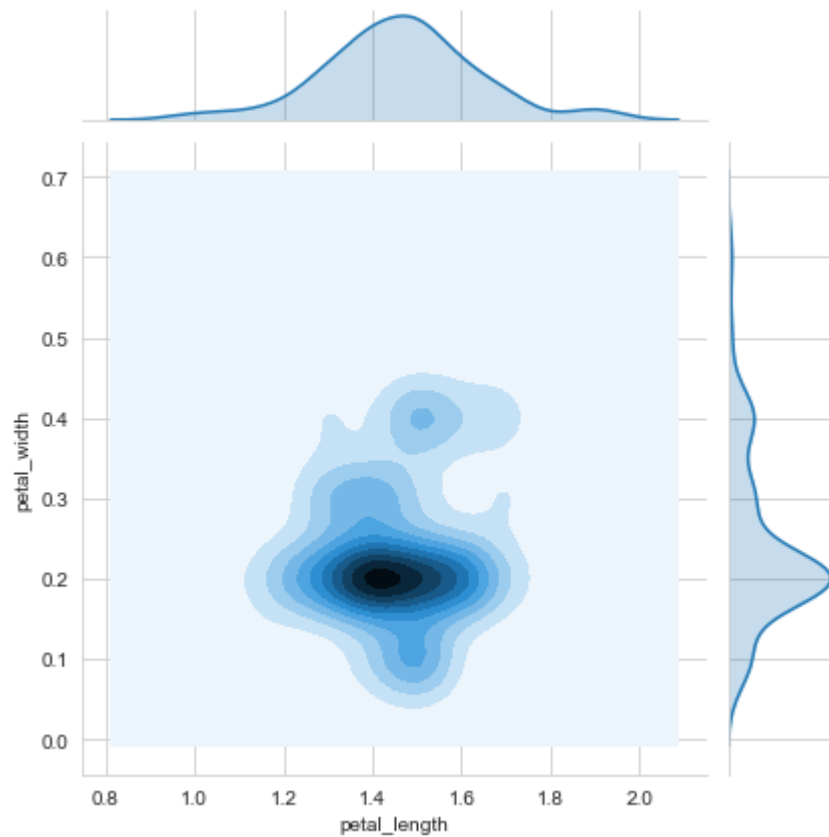


# violin plots

In [63]:
```
# 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()
```

In [64]:
```python
#2D Density plot, contors-plot
sns.jointplot(x="petal_length", y="petal_width", data=iris_setosa, kind="kde");
plt.show();
```



# Modeling with scikit learn

In [76]: 
```python
X = iris.drop([ 'species'], axis=1)
y = iris['species']
print(X.head())
print("_"*100)

print(X.shape)
print("_"*100)
print(y.head())
print("_"*100)
print(y.shape)
```

```
   sepal_length  sepal_width  petal_length  petal_width
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
_____
(150, 4)
_____
0    setosa
1    setosa
2    setosa
3    setosa
4    setosa
Name: species, dtype: object
_____
(150,)
```
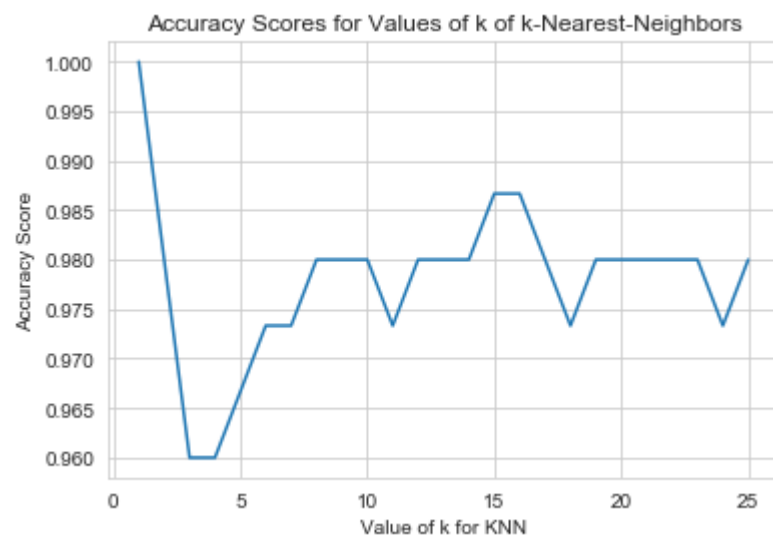
## Initialization the knn model

In [66]:
```python
k_range = list(range(1,26))
scores = []
for k in k_range:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X, y)
    y_pred = knn.predict(X)
    scores.append(metrics.accuracy_score(y, y_pred))

plt.plot(k_range, scores)
plt.xlabel('Value of k for KNN')
plt.ylabel('Accuracy Score')
plt.title('Accuracy Scores for Values of k of k-Nearest-Neighbors')
plt.show()
```



# Initialization Logistic Regression

```
In [67]:  logreg = LogisticRegression()
          logreg.fit(X, y)
          y_pred = logreg.predict(X)
          print(metrics.accuracy_score(y, y_pred))
```

0.96

```
C:\Users\SUNNY\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solve
r will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
C:\Users\SUNNY\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:469: FutureWarning: Default multi
_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
  "this warning.", FutureWarning)
```

**By splitting the dataset pseudo-randomly into a two separate sets, we can train using one set and test using another.**

This ensures that we won't use the same observations in both sets. More flexible and faster than creating a model using all of the dataset for training.

## Spiting the data into train and test data

```
In [68]:  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=5)
          print(X_train.shape)
          print(y_train.shape)
          print(X_test.shape)
          print(y_test.shape)
```
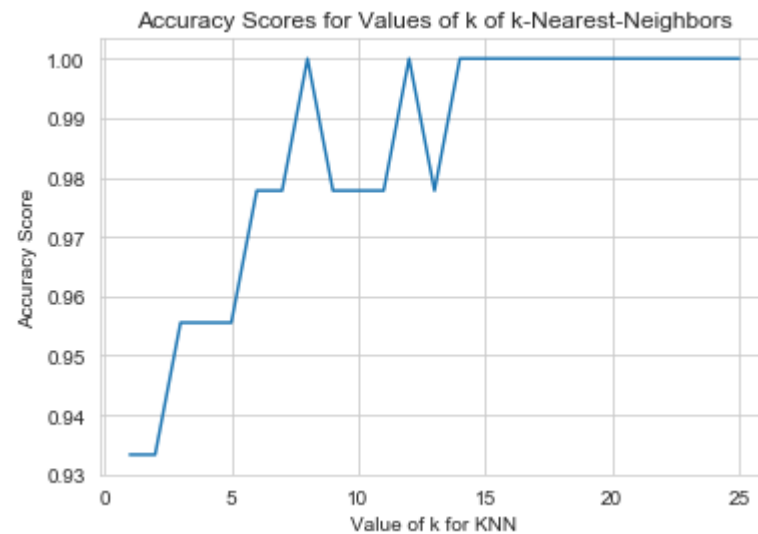
```
(105, 4)
(105,)
(45, 4)
(45,)
```

In [77]:
```python
##
#initialization the KNN MODEL

k_range = list(range(1,26))
scores = []
for k in k_range:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)
    scores.append(metrics.accuracy_score(y_test, y_pred))

plt.plot(k_range, scores)
plt.xlabel('Value of k for KNN')
plt.ylabel('Accuracy Score')
plt.title('Accuracy Scores for Values of k of k-Nearest-Neighbors')
plt.show()
```

Accuracy Scores for Values of k of k-Nearest-Neighbors

In [70]:
```
##  initialization Logistic regression


logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)
print(metrics.accuracy_score(y_test, y_pred))
```

0.9333333333333333

```
C:\Users\SUNNY\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solve
r will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
C:\Users\SUNNY\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:469: FutureWarning: Default multi
_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
  "this warning.", FutureWarning)
```

In [ ]:

In [ ]:

In [ ]:

In [ ]: