

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
sns.set_palette('husl')
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn import metrics
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split

data = pd.read_csv('../input/Iris.csv')
```

Preview of Data

- There are 150 observations with 4 features each (sepal length, sepal width, petal length, petal width).
- There are no null values, so we don't have to worry about that.
- There are 50 observations of each species (setosa, versicolor, virginica).

```
In [2]: data.head()
```

```
Out[2]:
```

	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

```
In [3]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
Id                150 non-null int64
SepalLengthCm    150 non-null float64
SepalWidthCm     150 non-null float64
PetalLengthCm    150 non-null float64
PetalWidthCm     150 non-null float64
Species          150 non-null object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.1+ KB
```

In [4]: `data.describe()`

Out[4]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	75.500000	5.843333	3.054000	3.758667	1.198667
std	43.445368	0.828066	0.433594	1.764420	0.763161
min	1.000000	4.300000	2.000000	1.000000	0.100000
25%	38.250000	5.100000	2.800000	1.600000	0.300000
50%	75.500000	5.800000	3.000000	4.350000	1.300000
75%	112.750000	6.400000	3.300000	5.100000	1.800000
max	150.000000	7.900000	4.400000	6.900000	2.500000

In [5]: `data['Species'].value_counts()`

Out[5]:

```
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
Name: Species, dtype: int64
```

Data Visualization

- After graphing the features in a pair plot, it is clear that the relationship between pairs of features of a iris-setosa (in pink) is distinctly different from those of the other two species.
- There is some overlap in the pairwise relationships of the other two species, iris-versicolor (brown) and iris-virginica (green).

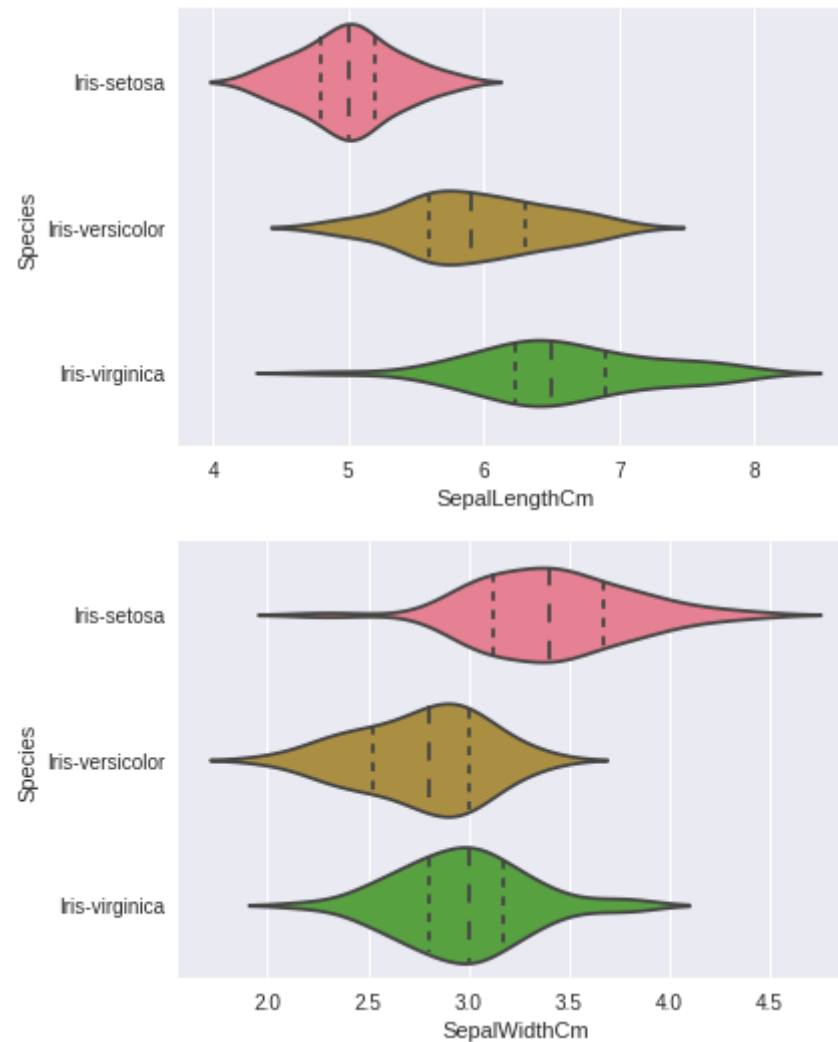
In [6]:

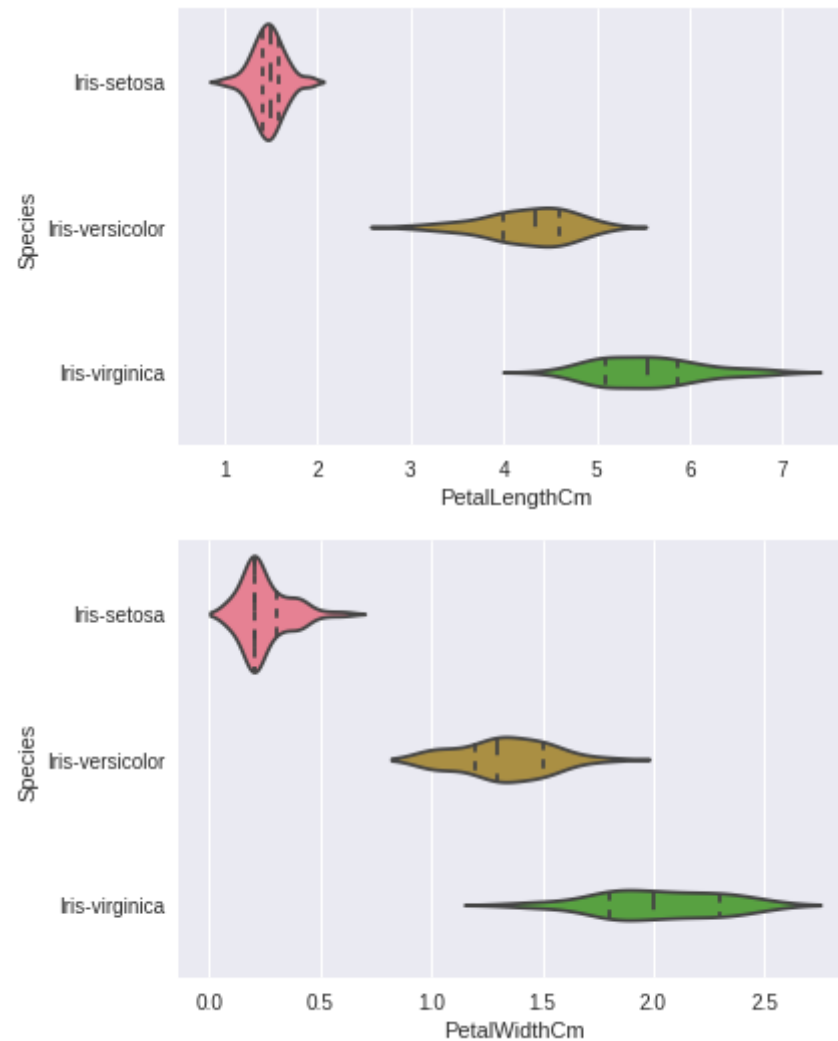
```
tmp = data.drop('Id', axis=1)
g = sns.pairplot(tmp, hue='Species', markers='+')
plt.show()
```



In [7]:

```
g = sns.violinplot(y='Species', x='SepalLengthCm', data=data, inner='quartile')  
plt.show()  
g = sns.violinplot(y='Species', x='SepalWidthCm', data=data, inner='quartile')  
plt.show()  
g = sns.violinplot(y='Species', x='PetalLengthCm', data=data, inner='quartile')  
plt.show()  
g = sns.violinplot(y='Species', x='PetalWidthCm', data=data, inner='quartile')  
plt.show()
```





Modeling with scikit-learn

```
In [8]: X = data.drop(['Id', 'Species'], axis=1)
y = data['Species']
# print(X.head())
print(X.shape)
# print(y.head())
print(y.shape)
```

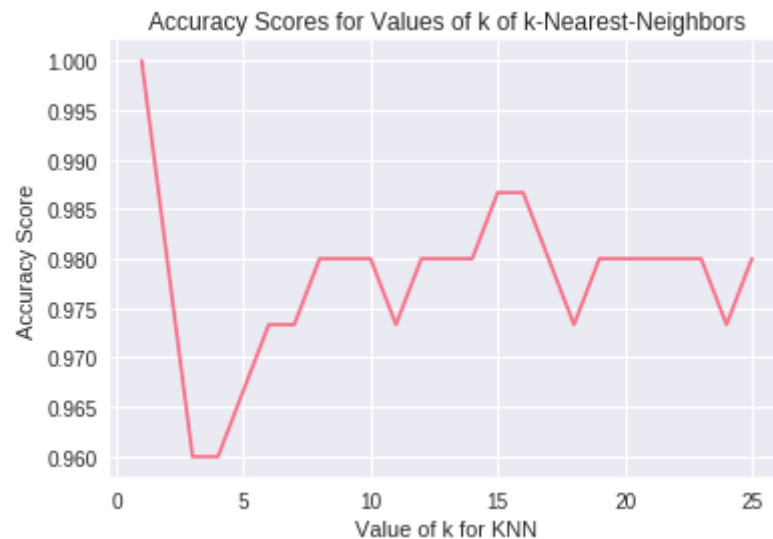
```
(150, 4)  
(150,)
```

Train and test on the same dataset

- This method is not suggested since the end goal is to predict iris species using a dataset the model has not seen before.
- There is also a risk of overfitting the training data.

In [9]:

```
# experimenting with different n values  
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()
```



```
In [10]: logreg = LogisticRegression()  
logreg.fit(X, y)  
y_pred = logreg.predict(X)  
print(metrics.accuracy_score(y, y_pred))
```

0.96

Split the dataset into a training set and a testing set

Advantages

- 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.

Disadvantages

- The accuracy scores for the testing set can vary depending on what observations are in the set.
- This disadvantage can be countered using k-fold cross-validation.

Notes

- The accuracy score of the models depends on the observations in the testing set, which is determined by the seed of the pseudo-random number generator (random_state parameter).
- As a model's complexity increases, the training accuracy (accuracy you get when you train and test the model on the same data) increases.
- If a model is too complex or not complex enough, the testing accuracy is lower.
- For KNN models, the value of k determines the level of complexity. A lower value of k means that the model is more complex.

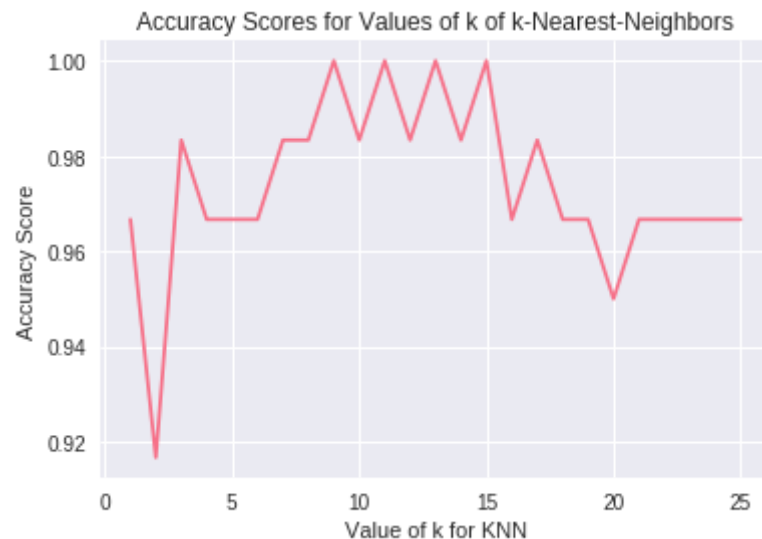
```
In [11]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=5)  
print(X_train.shape)  
print(y_train.shape)  
print(X_test.shape)  
print(y_test.shape)
```



```
(90, 4)  
(90,)   
(60, 4)  
(60,)
```

In [12]:

```
# experimenting with different n values  
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()
```



In [13]:

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

0.933333333333

Choosing KNN to Model Iris Species Prediction with $k = 12$

After seeing that a value of $k = 12$ is a pretty good number of neighbors for this model, I used it to fit the model for the entire dataset instead of just the training set.

```
In [14]: knn = KNeighborsClassifier(n_neighbors=12)
knn.fit(X, y)

# make a prediction for an example of an out-of-sample observation
knn.predict([[6, 3, 4, 2]])
```

```
Out[14]: array(['Iris-versicolor'], dtype=object)
```