

About our Dataset

Our Target

Implement a machine learning classification or regression model on the dataset. Classification is the task of separating items into its corresponding class.

```
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns
```

```
df = pd.read_csv('iris.data', header = None)
df
```

	0	1	2	3	4
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa

3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
...

145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

Giving the dataset its columns names

```
In [95]: df['class'].unique()
Out[95]: array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)
```

	sepal length	sepal width	petal length	petal width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.3	0.2	Iris-setosa
2	4.7	3.2	1.1	0.2	Iris-setosa
3	4.6	3.1	1.0	0.2	Iris-setosa
4	5.0	3.6	1.5	0.1	Iris-setosa
...
145	6.7	3.0	2.0	0.1	Iris-versicolor
146	6.3	2.5	2.3	0.3	Iris-versicolor
147	6.5	3.0	2.5	0.2	Iris-versicolor
148	6.2	3.4	2.3	0.2	Iris-versicolor
149	5.9	3.0	2.0	0.1	Iris-versicolor

150 rows × 5 columns

Checking for Null Values

```
In [29]: df.isna().sum()
```

```
Out[29]:
```

	sepal length	sepal width	petal length	petal width	class
count	0	0	0	0	0
dtype	int64				

```
In [30]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
 #   Column       Non-Null Count
 0   sepal length    150
 1   sepal width     150
 2   petal length    150
 3   petal width     150
 4   class          150
```

```
df.describe()
```

	sepal length	sepal width	petal length	petal width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

We do this to get the main attributes that should be used for classification.

```
In [16]: sns.pairplot(data = df, hue ='class')
```

Out[16]: <seaborn.axisgrid.PairGrid at 0x123d87fd7e0>

The figure is a 4x4 grid of plots representing the Iris dataset. The diagonal elements show histograms of individual attributes (sepal length, sepal width, petal length, petal width) for three classes: Iris-setosa (blue), Iris-versicolor (orange), and Iris-virginica (green). The off-diagonal elements show scatter plots of pairs of attributes (e.g., sepal length vs. sepal width, petal length vs. petal width) for each class. Each scatter plot includes marginal density curves along the axes.

```
#fig = plt.figure(figsize=(10,8))
fig, axes = plt.subplots(nrows = 1,ncols =
a = 0
for i in range(0,100,20):
```

```
tree.plot_tree(RFC.estimators_[i], feature_names =  
    filled=True, impurity=True, rounded=True , a  
  
axes[a].set_title('Estimators -: ' + str(i), fontsize=12)
```

```
fig.savefig('rf_5trees.png')
```

```
graph TD; A["petal length <= 2.37  
gini = 0.854  
impurity = 0.001  
value = [16, 47, 29]  
class = 'Iris-versicolor'] --> B["petal width <= 0.7  
gini = 0.66  
impurity = 0.001  
value = [1, 1, 1]  
class = 'Iris-versicolor']
```

The decision tree diagram illustrates the following splits:

- Root Node:** value = {0, 13, 2} (samples = 10, 0.515, 0.485)
- Left Child (value <= 2):** samples = 2, class = 0-10, impurity = 0.0
- Right Child (value > 2):** samples = 8, class = 11-20, impurity = 0.485
- Left Child of Right Node (value <= 4, 10):** samples = 5, class = 11-15, impurity = 0.0
- Right Child of Right Node (value > 4, 10):** samples = 3, class = 16-20, impurity = 0.485
- Left Child of Left Child (width <= 3.05):** samples = 2, class = 0-10, impurity = 0.0
- Right Child of Left Child (width > 3.05):** samples = 5, class = 11-15, impurity = 0.0
- Left Child of Right Child (width <= 4.95):** samples = 1, class = 0-10, impurity = 0.0
- Right Child of Right Child (width > 4.95):** samples = 1, class = 11-20, impurity = 0.0

```

graph TD
    Root["age < 10.0; samples = 1; is_independent"] --> S1["age = [0, 4]; samples = 2; is_independent"]
    Root --> S2["age = [0, 4]; samples = 2; is_independent"]
    S1 --> L1["age = 0.0; samples = 1; is_independent"]
    S1 --> L2["age = 0.5; samples = 1; is_independent"]
    S2 --> L3["age = 0.0; samples = 1; is_independent"]
    S2 --> L4["age = 0.5; samples = 1; is_independent"]
  
```

The decision tree diagram illustrates splits based on the 'age' feature and the number of 'samples'. The root node splits into two branches: one for 'age < 10.0' and 'samples = 1' (labeled 'is_independent'), and another for 'age >= 10.0' and 'samples >= 2' (also labeled 'is_independent'). Each branch further splits on 'age' into two child nodes. The first child node for the root splits into 'age = 0.0' and 'age = 0.5' (both labeled 'is_independent'). The second child node for the root splits into 'age = 0.0' and 'age = 0.5' (both labeled 'is_independent'). All leaf nodes are labeled as 'is_independent'.

Predicting the Xtesting values

```
from sklearn.metrics import
```

```
cm
```

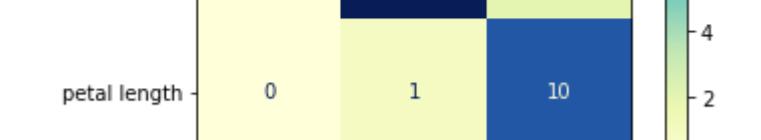
```
Out[57]: array([[12,  0,  0],
                  [ 0, 13,  2],
                  [ 0,  1, 10]], dtype=int64)
```

A 2x2 confusion matrix heatmap. The x-axis is labeled 'Actual' and the y-axis is labeled 'Predicted'. The diagonal elements are dark blue, representing correct classifications. The off-diagonal elements are light yellow, representing misclassifications. The values are: Actual 0, Predicted 0: 12; Actual 0, Predicted 1: 2; Actual 1, Predicted 0: 2; Actual 1, Predicted 1: 12.

		Actual 0	Actual 1
Predicted 0	12	2	
Predicted 1	2	12	

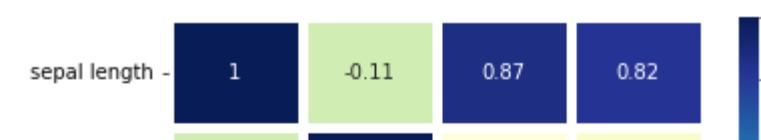
The figure is a heatmap illustrating the distribution of sepal width across three clusters. The x-axis represents sepal width values (0, 13, 2) and the y-axis represents the cluster index (0, 1, 2). The color scale ranges from dark blue (low value) to light yellow (high value).

Cluster	Sepal Width 0	Sepal Width 13	Sepal Width 2
0	High (Yellow)	Low (Dark Blue)	Medium (Light Green)
1	Medium (Light Green)	Very Low (Dark Blue)	Very Low (Dark Blue)
2	Medium (Light Green)	Very Low (Dark Blue)	Very Low (Dark Blue)



```
In [77]: from sklearn.metrics import accuracy_score
```

```
In [28]: sns.heatmap(df.corr(), cmap = 'YlGnBu' , linewidths=1)  
Out[28]: <AxesSubplot:>
```



sepal width -	-0.11	1	-0.42	-0.36
petal length -	0.87	-0.42	1	0.96

```
x = df.iloc[:, :-1]
```