

EDA on the Iris dataset

Data Scientist : Simon Pierre Charbel

```
In [1]: 1 # Import Libraries and packages
2 import pandas as pd
3 from sklearn.datasets import load_iris
4
5 import seaborn as sns
6 import matplotlib.pyplot as plt
```

1. Loading and pre-processing the data

```
In [2]: 1 # Load the Iris dataset
2 iris = load_iris()
3 print(f"Type of Scikit-Learn object returned : {type(iris)}\n")
4
5 # Create a pandas DataFrame
6 iris_df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
7
8 # Add the target column to the DataFrame
9 iris_df['target'] = iris.target
10 # Add the target names as a new column
11 iris_df['species'] = iris.target_names[iris.target]
12
13 # Print shape and the first few rows of the dataframe
14 print(f"DataFrame Shape: {iris_df.shape[0]} rows and {iris_df.shape[1]}")
15 iris_df.head()
```

Type of Scikit-Learn object returned : <class 'sklearn.utils._bunch.Bunch'>

DataFrame Shape: 150 rows and 6 columns (shape).

Out[2]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target	species
0	5.1	3.5	1.4	0.2	0	setosa
1	4.9	3.0	1.4	0.2	0	setosa
2	4.7	3.2	1.3	0.2	0	setosa
3	4.6	3.1	1.5	0.2	0	setosa
4	5.0	3.6	1.4	0.2	0	setosa

```
In [3]: 1 # Checking for missing values
        2 iris_df.isnull().sum()
```

```
Out[3]: sepal length (cm)    0
        sepal width (cm)     0
        petal length (cm)    0
        petal width (cm)     0
        target               0
        species              0
        dtype: int64
```

```
In [4]: 1 # Check for duplicate rows
        2 duplicates = iris_df.duplicated(keep=False)
        3 # Display all rows that are duplicates
        4 duplicate_rows = iris_df[duplicates]
        5 #print( f" Duplicate rows : \n {duplicate_rows}")
        6 print( f" Duplicate rows count : {duplicate_rows.shape[0]}")
        7
        8 # Drop the Duplicates
        9 iris_df.drop_duplicates(inplace=True)
```

```
Duplicate rows count : 2
```

```
In [5]: 1 # Print the shape of the DataFrame
2 print(f"New DataFrame Shape: {iris_df.shape[0]} rows and {iris_df.shape[1]} columns")
3
4 # Print the column types
5 print(f"\nColumns types: \n")
6 print(iris_df.info())
7
8 # Print a statistical summary of the dataset
9 print(f"\n\n A statistical summary of the dataset:")
10 iris_df.describe()
```

New DataFrame Shape: 149 rows and 6 columns (shape).

Columns types:

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 149 entries, 0 to 149
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  -
0   sepal length (cm)      149 non-null   float64
1   sepal width (cm)       149 non-null   float64
2   petal length (cm)      149 non-null   float64
3   petal width (cm)       149 non-null   float64
4   target                 149 non-null   int32
5   species                149 non-null   object
dtypes: float64(4), int32(1), object(1)
memory usage: 7.6+ KB
None
```

A statistical summary of the dataset:

Out[5]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
count	149.000000	149.000000	149.000000	149.000000	149.000000
mean	5.843624	3.059732	3.748993	1.194631	0.993289
std	0.830851	0.436342	1.767791	0.762622	0.817847
min	4.300000	2.000000	1.000000	0.100000	0.000000
25%	5.100000	2.800000	1.600000	0.300000	0.000000
50%	5.800000	3.000000	4.300000	1.300000	1.000000
75%	6.400000	3.300000	5.100000	1.800000	2.000000
max	7.900000	4.400000	6.900000	2.500000	2.000000

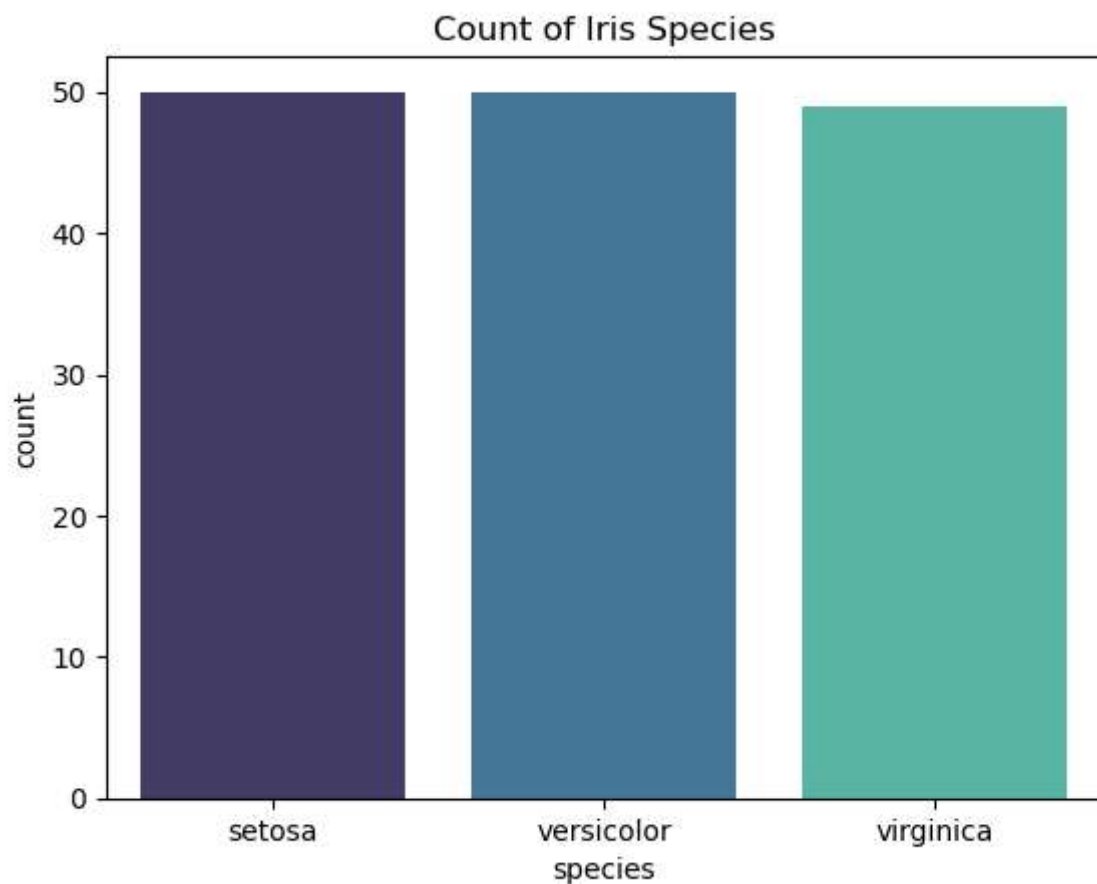
```
In [6]: 1 # Check how many species we have and their counts
        2 iris_df.value_counts("species")
```

```
Out[6]: species
setosa      50
versicolor  50
virginica   49
dtype: int64
```

2. Data Visualization

```
In [7]: 1 # Visualizing the target column (Countplot)
        2 sns.countplot(data=iris_df, x='species', palette='mako')
        3 plt.title("Count of Iris Species", fontsize=12)
        4 plt.figure(figsize=(4, 4))
```

```
Out[7]: <Figure size 400x400 with 0 Axes>
```

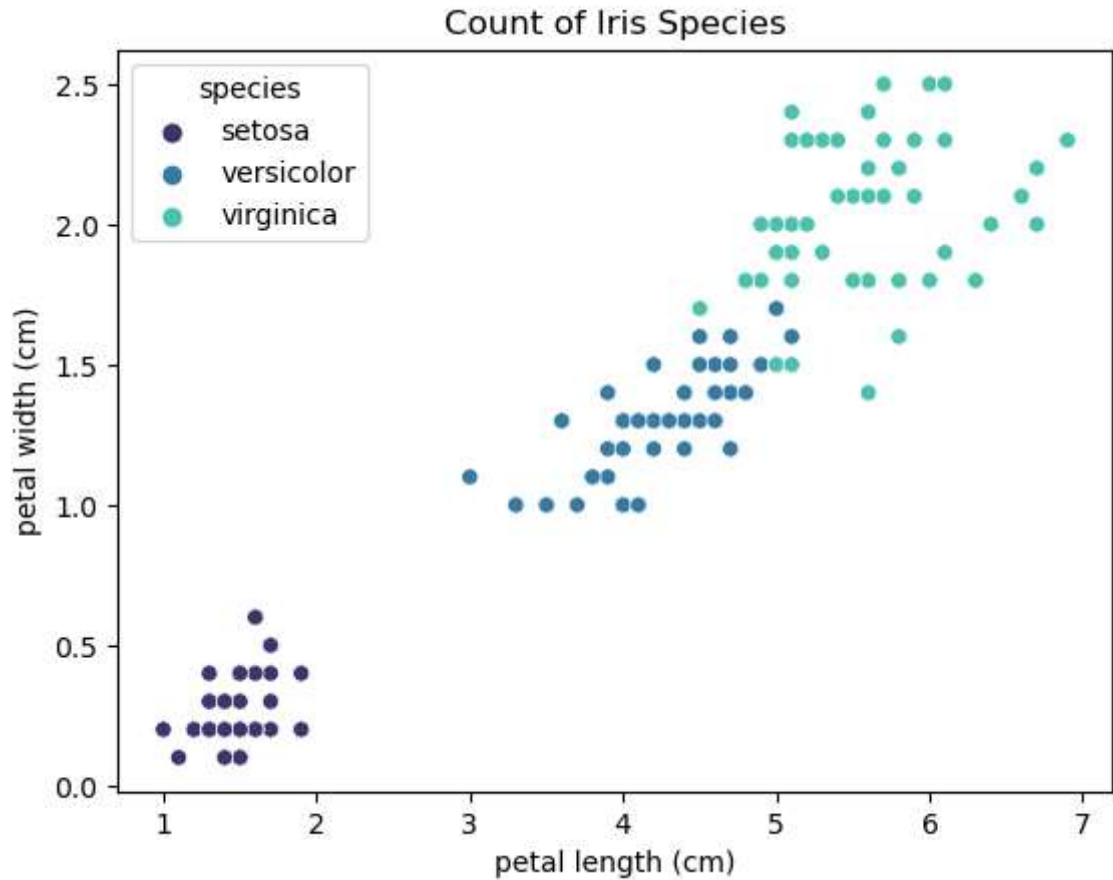


```
<Figure size 400x400 with 0 Axes>
```

- Explore the correlations between sepal length and sepal width, as well as between petal length and petal width.

```
In [8]: 1 # Petal
2 sns.scatterplot(data=iris_df, x='petal length (cm)',y='petal width (cm)
3             hue='species', palette='mako')
4 plt.title("Count of Iris Species", fontsize=12)
5 plt.figure(figsize=(4, 4))
```

Out[8]: <Figure size 400x400 with 0 Axes>

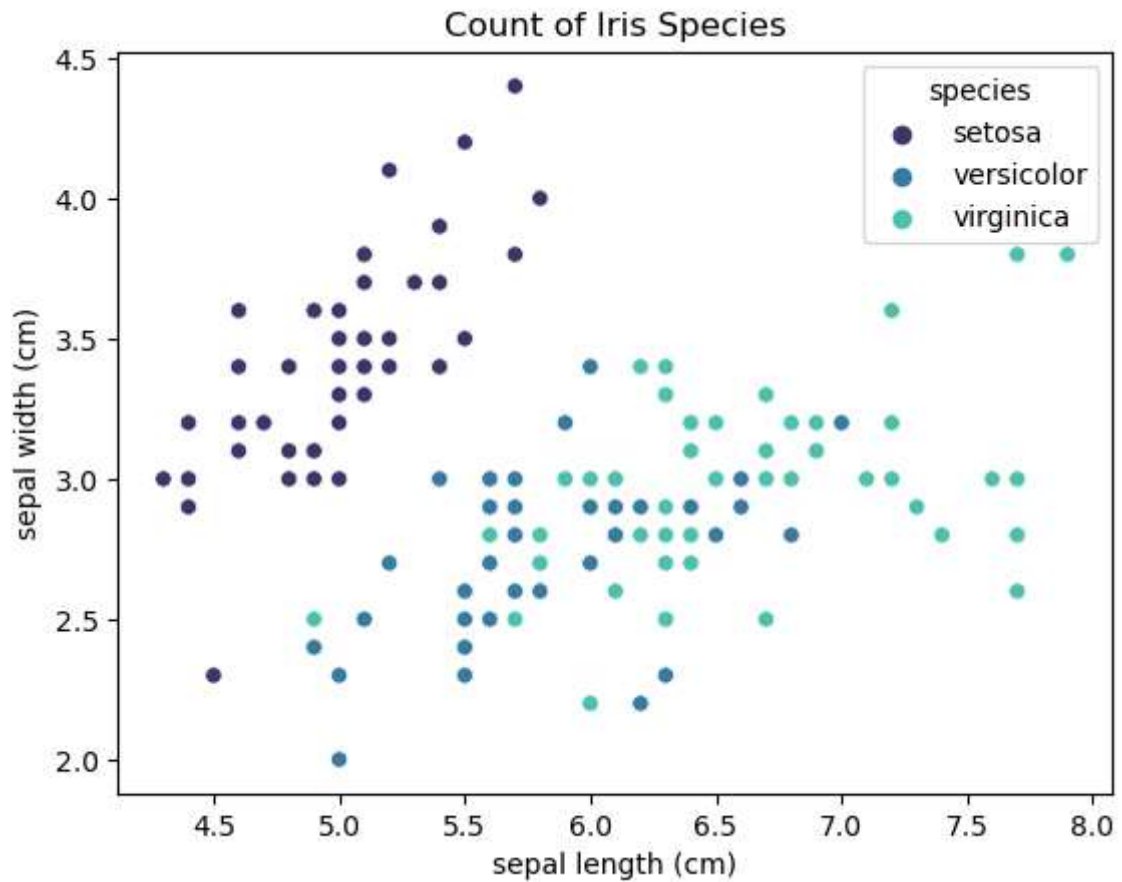


<Figure size 400x400 with 0 Axes>

- The Setosa species exhibits comparatively smaller petal lengths and widths.
- The Versicolor species falls in between the other two species regarding petal length and width.
- The Virginica species boasts the largest petal lengths and widths among them.

```
In [9]: 1 # Septal
2 sns.scatterplot(data=iris_df, x='sepal length (cm)',y='sepal width (cm)
3             hue='species', palette='mako')
4 plt.title("Count of Iris Species", fontsize=12)
5 plt.figure(figsize=(4, 4))
```

Out[9]: <Figure size 400x400 with 0 Axes>



<Figure size 400x400 with 0 Axes>

- The Setosa species showcases smaller sepal lengths but wider sepal widths.
- The Versicolor species occupies an intermediate position between the other two species concerning sepal length and width.
- The Virginica species boasts longer sepal lengths but narrower sepal widths.

```
In [158]: 1 file_path = r'C:\Users\simon\AppData\Roaming\JetBrains\PyCharmCE2022.2\
2
3 # Save the DataFrame to a CSV file with UTF-8 encoding
4 iris_df.to_csv(file_path, encoding='utf-8', index=False)
```