# How to interpret pair plot?



Pairwise Relationships in Iris Dataset (Numerical Variables) by Species
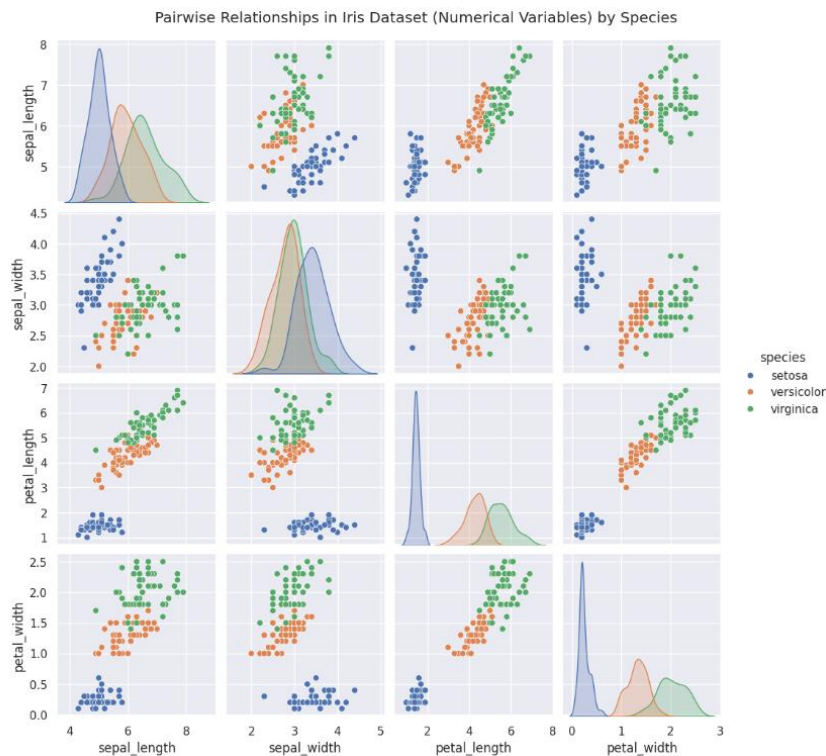
## A. Understanding the Components of a Pair Plot:

A pair plot creates a grid of plots where:

- **Diagonal:** The diagonal plots show the univariate distribution of each numerical variable. Typically, this is represented as a histogram or a kernel density estimate (KDE). In this image, KDE plots are used, with each species having a different color.
- **Off-Diagonal:** The off-diagonal plots show the bivariate relationship between each pair of numerical variables. These are typically scatter plots. In this image, the points are colored by the species.

The variables displayed are: sepal_length, sepal_width, petal_length, and petal_width. The species are setosa (blue), versicolor (orange), and virginica (green).

1. **Diagonal Plots (Univariate Distributions by Species):**

   - **sepal_length (Top-Left):** The KDE plot shows the distribution of sepal length for each species. setosa has a lower sepal length compared to versicolor and virginica, which have overlapping but distinct distributions.
   - **sepal_width (Middle-Left):** The KDE plot shows the distribution of sepal width. setosa generally has a larger sepal width compared to versicolor and virginica, which have more overlap.
   - **petal_length (Middle-Right):** The KDE plot shows the distribution of petal length. The species are clearly separated here, with setosa having the smallest petal length, followed by versicolor, and then virginica with the largest.
   - **petal_width (Bottom-Right):** Similar to petal length, the KDE plot for petal width shows good separation between the species, with setosa having the smallest, followed by versicolor, and then virginica with the largest.

2. **Off-Diagonal Plots (Bivariate Relationships by Species):**

   - **sepal_length vs. sepal_width (Top-Middle-Left):** The scatter plot shows the relationship between sepal length and sepal width for each species. setosa tends to have shorter but wider sepals. versicolor and virginica have more overlap, with virginica generally having longer and slightly narrower sepals than versicolor.
   - **sepal_length vs. petal_length (Top-Middle-Right):** There's a clear positive correlation between sepal length and petal length, especially for versicolor and virginica. setosa has short sepals and very short petals, forming a distinct cluster.
   - **sepal_length vs. petal_width (Top-Right):** Similar to petal length, there's a positive correlation between sepal length and petal width, with setosa forming a separate cluster with small petal widths.
   - **sepal_width vs. petal_length (Middle-Left-Bottom):** There's a less clear correlation here. setosa has wider sepals and short petals. versicolor and virginica show some positive trend but with more spread.

- **sepal_width vs. petal_width (Middle-Right-Bottom):** Similar to sepal width vs. petal length, the correlation is not very strong. setosa has wider sepals and narrow petals. versicolor and virginica show some positive trend.
- **petal_length vs. petal_width (Bottom-Left):** There's a strong positive correlation between petal length and petal width. The species form distinct clusters, especially setosa with small petals and virginica with large petals. versicolor lies in between.

## C. Overall Interpretation:

The pair plot provides a comprehensive view of the relationships between the four numerical features in the Iris dataset, broken down by species. It reveals:

- **Separability of Species:** The features, particularly petal length and petal width, are very effective in distinguishing the three species. setosa forms a clearly separated cluster in many of the bivariate plots. versicolor and virginica have some overlap but can often be distinguished, especially using petal dimensions.
- **Correlations Between Features:** There are noticeable correlations between certain features, such as the positive correlation between petal length and petal width, and between sepal length and petal length (for versicolor and virginica).
- **Univariate Distributions:** The diagonal plots show the range and shape of the distribution for each feature within each species.

## Pair plots are most valuable when you want to:

- **Explore the pairwise relationships between multiple numerical variables in a dataset.** They provide a quick visual overview of potential correlations and distributions.
- **Identify potential variables that might be useful for separating different groups or classes in your data.** In this case, the clear separation of species based on petal length and width is evident.
- **Detect potential multicollinearity between numerical features.** Strong correlations in the off-diagonal plots can indicate multicollinearity, which might be a concern for some modeling techniques.

- **Get a sense of the overall structure and patterns in a multivariate dataset.** They allow you to see both the individual distributions and the joint relationships.
- **Perform initial exploratory data analysis (EDA) before applying more complex modeling techniques.** They can help generate hypotheses about the relationships between variables.
- **Visualize the relationships in a classification problem where you want to see how well the features can discriminate between different classes.** Coloring the points by class (as done by species here) is very informative.

In summary, pair plots are an excellent tool for gaining initial insights into the multivariate relationships within a dataset of numerical variables, especially when you also have a categorical variable (like 'species' here) to color-code the observations. They provide a comprehensive visual summary of both univariate and bivariate distributions.

N.B: Pair plot can be used for both Bivariate and Multivariate Analysis