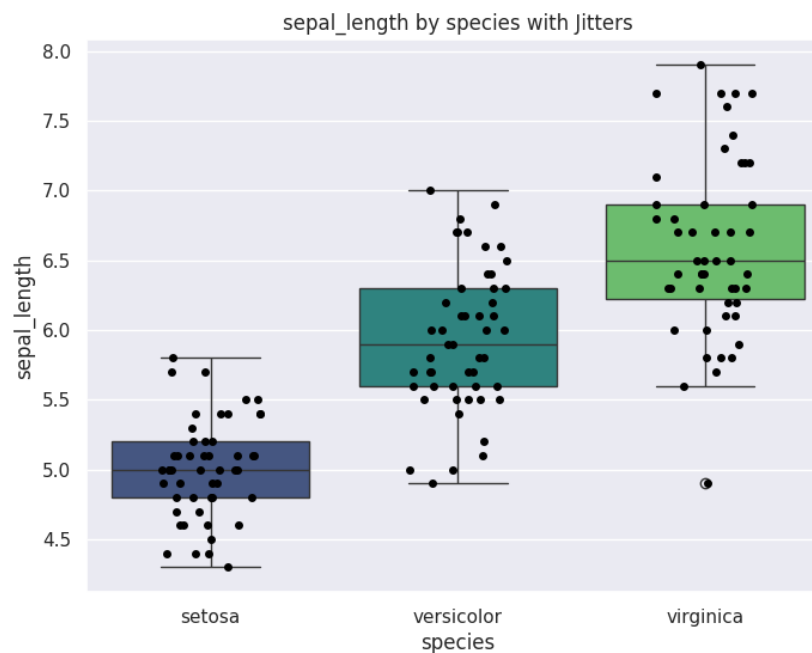


## How to interpret Boxplot with Jitters for bivariate analysis



### A. Understanding the Components of a Box Plot with Jitters:

- **Horizontal Axis (X-axis):** Represents the categorical variable "species," with three categories: "setosa," "versicolor," and "virginica."
- **Vertical Axis (Y-axis):** Represents the numerical variable "sepal\_length."
- **Boxes:** Each box summarizes the distribution of sepal length for a specific species:
  - The **bottom line** of the box indicates the first quartile (Q1), where 25% of the data falls below.
  - The **top line** of the box indicates the third quartile (Q3), where 75% of the data falls below.
  - The **line inside the box** indicates the median (Q2), the middle value of the data.
  - The **height of the box** represents the interquartile range (IQR =  $Q3 - Q1$ ), which contains the middle 50% of the data.

- **Whiskers:** The lines extending from the top and bottom of the box are the whiskers. They typically extend to 1.5 times the IQR from the quartiles, showing the spread of the majority of the data. Data points outside the whiskers are considered potential outliers.
- **Outliers:** Individual points plotted outside the whiskers represent values that are unusually high or low compared to the rest of the data for that species.
- **Jitters (Black Dots):** The black dots are the individual data points for each species, slightly scattered horizontally (jittered) to prevent overlap and show the density of observations at different sepal lengths within each species.

## B. Interpreting the Sepal Length Distribution by Species:

By examining the box plots and the jittered points, we can compare the sepal length distributions across the three species:

- **Setosa:**
  - Has the shortest median sepal length (around 5.0 cm).
  - The box is relatively short, indicating less variability in sepal length compared to the other species.
  - The jittered points are clustered tightly, further showing low variability.
  - There are some points at the lower end that might be considered lower outliers.
- **Versicolor:**
  - Has a median sepal length around 5.9 cm, which is longer than setosa but shorter than virginica.
  - The box is taller than that of setosa, indicating more variability in sepal length.
  - The jittered points are more spread out along the sepal length range.
  - There are no obvious high or low outliers.

- **Virginica:**
  - Has the longest median sepal length (around 6.5-6.6 cm).
  - The box is similar in height to that of versicolor, indicating comparable variability.
  - The jittered points are also spread out, showing a range of sepal lengths.
  - There is one potential low outlier.

### C. Overall Interpretation:

The box plot with jitters clearly shows that there are differences in sepal length among the three species of Iris. Setosa generally has shorter sepals with less variability, while versicolor has intermediate sepal lengths with more variability, and virginica has the longest sepal lengths with a similar degree of variability to versicolor. The jittered points provide a more detailed view of the distribution and density of sepal length values for each species, complementing the summary statistics presented by the box plots.

### Box plots with jitters are particularly useful when you want to:

- **Compare the distribution of a numerical variable across different categories of a categorical variable.** The boxes provide a clear summary of the central tendency, spread, and skewness of the numerical data for each category.
- **Visualize the individual data points in addition to the summary statistics.** The jitters allow you to see the actual data distribution and identify clusters or gaps that might be hidden by the box plot alone.
- **Identify potential outliers within each category.** The points outside the whiskers are easily noticeable.
- **Get a sense of the density of the data points within each category.** Areas with more overlapping jittered points indicate higher density.
- **Compare the variability of the numerical variable across different categories.** The height of the boxes and the spread of the jitters provide insights into the dispersion of the data.

- **Present a concise yet informative visualization that combines summary statistics with the raw data.** This can be particularly helpful for communicating findings to both technical and non-technical audiences.

In summary, box plots with jitters offer a powerful way to compare the distributions of a numerical variable across different categories by combining the strengths of box plots (summary statistics) with the detailed information provided by individual data points (jitters). They are excellent for exploring and presenting differences in central tendency, spread, and the presence of outliers.