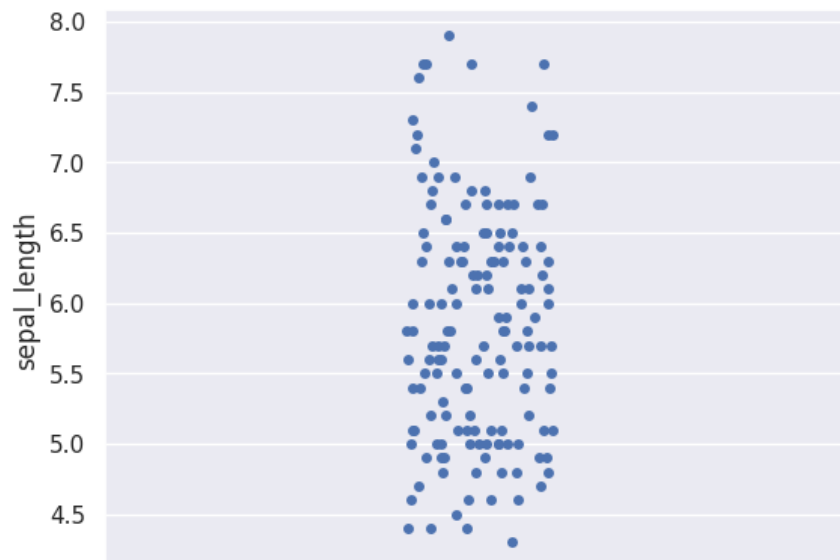


## How to Interpret Strip plot?



### A. Interpretation of the Strip Plot Components:

- **Vertical Axis (Y-axis):** Represents the range of values for the numerical variable, "sepal\_length," spanning from approximately 4.3 to 8.0. Each dot's vertical position corresponds to its sepal length value.
- **Horizontal Axis (X-axis):** This axis represents the single category being analyzed (in this univariate case, there's no grouping). The points are spread out horizontally (jittered) to prevent overlap and make it easier to see the density of data points at different sepal length values.
- **Dots (Individual Data Points):** Each dot represents a single observation in the dataset. Its vertical position indicates the value of "sepal\_length" for that observation.

### B. Interpreting the "sepal\_length" Distribution:

- **Central Tendency:** While a strip plot doesn't explicitly show summary statistics like the mean or median, you can visually estimate the central tendency by observing where the dots are most concentrated along the y-axis. In this plot, the highest density of points appears to be around the 5.5 to 6.5 range.

- **Spread (Variability):** The vertical range covered by the dots indicates the overall spread or variability of the sepal length measurements. Here, the data ranges from approximately 4.3 to 7.9.
- **Shape:** By looking at the density of the dots along the y-axis, you can get a sense of the distribution's shape:
  - Areas with more densely packed dots indicate higher frequencies of those sepal length values.
  - Areas with fewer dots indicate lower frequencies.
  - In this plot, the distribution seems somewhat uneven, with denser regions around 5.5-6.5 and fewer points at the extremes. It doesn't immediately suggest a perfectly normal or strongly skewed distribution, but there might be some subtle patterns.
- **Outliers:** Individual dots that are isolated and far from the main cluster of points can be identified as potential outliers. There appear to be a few points at the higher end (above 7.5) and possibly one at the lower end (around 4.3) that are somewhat separated.

**Strip plots are particularly useful for univariate analysis in the following scenarios:**

- **Small to Moderately Sized Datasets:** When the number of data points is not too large, a strip plot can effectively show each individual observation without significant overlap, allowing you to see the raw data.
- **Highlighting Individual Data Points:** If it's important to visualize every single data point and understand its exact value, a strip plot is a direct way to do this.
- **Detecting Gaps and Clusters:** The spacing between the jittered points can reveal gaps in the distribution (ranges with no or few observations) or clusters of data points at specific values.
- **Comparing Distributions Across Categories (Bivariate/Multivariate):** Strip plots are very effective when used to compare the distribution of a numerical variable across different categories of a categorical variable. By plotting the strip plot for each category side-by-side (along the x-

axis), you can easily see differences in central tendency, spread, and the distribution of individual points between the groups.

- **As a Complement to Summary Statistics or Other Visualizations:** A strip plot can be a valuable addition to box plots, violin plots, or histograms, providing the underlying raw data points that those summaries represent.

#### **In contrast to histograms, box plots, and violin plots:**

- **Histograms:** Group data into bins, losing the individual data point information. Strip plots retain each point.
- **Box Plots:** Provide a summary of the distribution (quartiles, median, outliers) but don't show the shape or individual points within those summaries. Strip plots show all the points.
- **Violin Plots:** Show the estimated probability density and summary statistics but not the individual data points as clearly as a strip plot with jitter.

In summary, strip plots are the best choice when you want to visualize each individual data point of a numerical variable, especially for small to moderate datasets, and when you want to see the raw data alongside or instead of summary statistics or binned representations. They are also excellent for comparing distributions across categories.