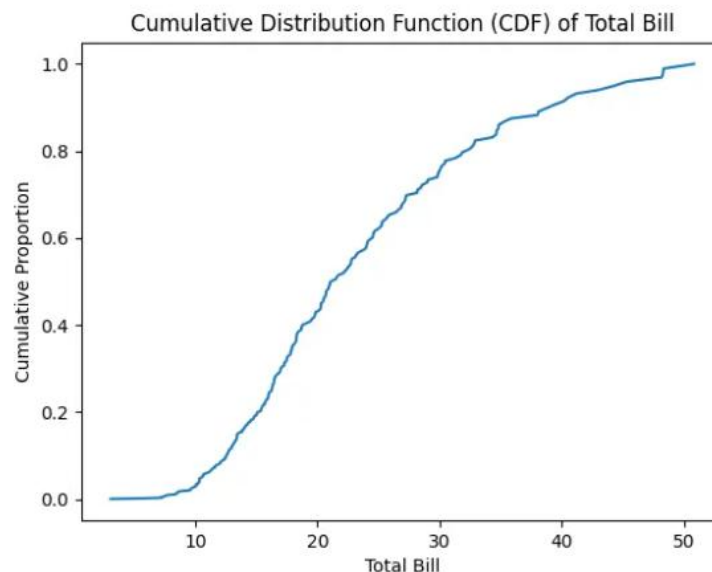


## How to interpret Cumulative Density Plot?



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

- **Horizontal Axis (X-axis):** Represents the range of values for the numerical variable, "Total Bill," spanning from approximately 0 to 55.
- **Vertical Axis (Y-axis):** Represents the **Cumulative Proportion** (or cumulative probability), ranging from 0 to 1. For any given value on the x-axis, the corresponding value on the y-axis indicates the proportion of data points that have a value *less than or equal to* that x-axis value.
- **The Line (CDF Curve):** The blue line shows the cumulative proportion as the "Total Bill" value increases.

### Interpreting the Curve:

- **Starting Point:** The curve starts near 0 on the y-axis for small values of "Total Bill" (around 5). This indicates that very few data points have a total bill less than or equal to 5.
- **Rising Trend:** As the "Total Bill" increases, the cumulative proportion also increases. This is because more data points fall below or equal to the increasing "Total Bill" values.

- **Steep Slopes:** Steeper sections of the curve indicate ranges where a larger proportion of the data points are concentrated. For example, the curve rises relatively steeply between a "Total Bill" of around 10 and 20, suggesting that a significant portion of the bills fall within this range.
- **Gradual Slopes (Plateaus):** Shallower sections of the curve indicate ranges where fewer data points are present. For instance, the curve rises more gradually for "Total Bill" values above 40, suggesting fewer bills in that higher range. The near-horizontal sections at the beginning and end indicate very few or no data points below or above those ranges, respectively (after accounting for potential outliers not explicitly shown).
- **Specific Points:**
  - To find the proportion of bills less than or equal to a specific value (e.g., 25), you would go to 25 on the x-axis and find the corresponding value on the y-axis. In this case, it's around 0.6. This means approximately 60% of the total bills are \$25 or less.
  - Similarly, to find the value below which a certain proportion of bills fall (e.g., the median, which corresponds to a cumulative proportion of 0.5), you would go to 0.5 on the y-axis and find the corresponding value on the x-axis. Here, the 50th percentile (median) is around 20.
- **Reaching 1:** The curve eventually reaches 1 on the y-axis as the "Total Bill" increases to its maximum value. This signifies that 100% of the data points have a "Total Bill" less than or equal to the highest observed value.

### Overall Interpretation of the "Total Bill" Distribution:

The CDF plot shows that:

- A significant proportion of the total bills are relatively low, concentrated in the range of \$10 to \$20.
- The median total bill is around \$20.

- The distribution is right-skewed. The steep rise in the lower values and the gradual increase in the higher values indicate that more bills are on the lower end, with a longer tail extending towards higher bill amounts.
- Very high total bills (above \$40-\$45) are less common.

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

- **Understanding the Proportion of Data Below a Certain Value:**  
The primary strength of a CDF is directly showing the cumulative proportion or probability. If you need to quickly answer questions like "What percentage of bills are less than \$30?" or "What is the 75th percentile of the total bill?", the CDF is the most direct visualization.
- **Comparing Distributions:** You can plot CDFs of multiple datasets on the same axes to easily compare their distributions. Differences in the steepness and position of the curves can reveal differences in central tendency, spread, and skewness. For example, if one CDF is consistently to the left of another, it indicates that the first dataset tends to have smaller values.
- **Determining Percentiles and Quantiles:** CDFs make it straightforward to find specific percentiles (e.g., 25th, 50th, 75th) or any other quantile of the distribution. You simply find the desired cumulative proportion on the y-axis and read the corresponding value on the x-axis.
- **Assessing the Empirical Distribution:** The CDF provides a direct visualization of the empirical cumulative distribution function of the data, without making assumptions about the underlying theoretical distribution.
- **When the Order of Data Matters:** CDFs inherently consider the ordered nature of the data, showing how the cumulative proportion builds as the variable's value increases.
- **In Reliability Analysis and Survival Analysis:** CDFs (often called failure functions in these contexts) are fundamental for understanding the probability of an event (e.g., failure of a component, death) occurring by a certain time.

### **In contrast to histograms and density plots:**

- **CDFs focus on the cumulative proportion, while histograms and density plots show the frequency or probability density at specific values.**
- **CDFs don't require binning (like histograms) or smoothing (like density plots), providing a direct representation of the empirical distribution.**
- **CDFs are particularly useful for comparing distributions and finding percentiles, tasks that might be less direct with histograms or density plots.**

In summary, CDF plots are the best choice when you need to understand the proportion of data below certain thresholds, compare distributions based on cumulative probabilities, and easily determine percentiles and quantiles of a numerical variable.