

Learning to Calculate and Plot Cumulative Distribution Functions (CDFs) in Python

Authored by
Mohammed loot

November 3, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning to Calculate and Plot Cumulative Distribution Functions (CDFs) in Python*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=9310>

The [Cumulative Distribution Function \(CDF\)](#) stands as a cornerstone in classical statistics, providing a comprehensive description of the probability distribution for a real-valued random variable. In the realm of modern data analysis and scientific computing, particularly when utilizing the [Python](#) ecosystem, the ability to accurately calculate and visualize the CDF is paramount for deciphering the inherent structure and distributional characteristics of any given dataset. It moves beyond simple summary statistics, offering a full picture of data spread.

For datasets derived from observed reality--known as empirical data--we primarily focus on computing the **Empirical CDF (eCDF)**. This process involves a systematic approach: first, sorting all observations in ascending order, and subsequently, assigning a cumulative probability rank to each data point. This non-parametric method requires no prior assumptions about the underlying theoretical distribution. The robust nature of the [NumPy](#) library in [Python](#) makes this calculation highly efficient, as demonstrated in the foundational code snippet below.

```
#sort data
```

```
x = np.sort(data)
```

```
#calculate CDF values
```

```
y = 1. * np.arange(len(data)) / (len(data) - 1)
```

```
#plot CDF
```

```
plt.plot(x, y)
```

Mastering this fundamental computational technique grants data scientists the capacity to visualize the true distribution of their data without relying on potentially flawed assumptions regarding a theoretical model. This methodology is indispensable in exploratory data analysis. The forthcoming sections will elaborate on the practical application of this core concept, providing step-by-step examples that differentiate between calculating the eCDF for raw, unstructured samples and leveraging specialized functions when the underlying theoretical distribution is already known.

Defining the Cumulative Distribution Function (CDF)

The [CDF](#), often denoted as $F(x)$, is formally defined as the probability that a random variable X will take on a value less than or equal to a specified value x . Unlike the **Probability Density Function (PDF)**, which illustrates the relative frequency or likelihood of a specific value occurring, the CDF is fundamentally cumulative. This means its range is strictly bounded between 0 and 1, providing a direct measure of accumulated probability. This cumulative nature makes the CDF an exceptionally intuitive tool for statistical inference, enabling immediate interpretation of percentiles and facilitating robust comparisons across disparate datasets.

In practical terms, the CDF answers the critical question: "What proportion or percentage of the

data falls below a given threshold value?" When visually inspecting a CDF plot, the slope provides immediate insight into data concentration. A significantly steep incline indicates a high density of data points clustered within that range, suggesting a mode or peak in the distribution. Conversely, a flatter slope signifies a sparse distribution where data points are widely spread. Because the CDF is mathematically guaranteed to be **monotonic** (always non-decreasing), it furnishes a stable and highly reliable representation of data spread, effectively mitigating the common visualization challenges, such as arbitrary binning choices, that frequently plague histograms or traditional PDF plots.

Crucially, for real-world statistical analysis, particularly when the precise theoretical form of the population distribution remains unknown, we rely heavily on the **Empirical CDF (eCDF)**. The eCDF is calculated directly from the observed sample data points. By sorting the values and determining the normalized rank of each observation, the eCDF serves as the best possible approximation of the true, underlying theoretical CDF based on the available data. This empirical calculation is vital for non-parametric statistical methods and underpins numerous modern statistical tests.

Computational Steps for Calculating the Empirical CDF

To calculate the eCDF manually within the [Python](#) environment, the process boils down to two distinct, yet equally important, steps: precise data sorting and subsequent rank normalization. The sorting phase, typically executed using the function `np.sort(data)` from the **NumPy** library, is non-negotiable. It ensures that the observations are processed sequentially in ascending order, which is the necessary prerequisite for accurately calculating the cumulative probability at each step.

The second, more nuanced step involves generating the cumulative probability values, which correspond to the y variable in our initial scripting examples. We utilize the [NumPy](#) function `np.arange(len(data))`, which systematically creates an index array spanning from 0 up to $N-1$, where N is the total count of observations. This array of indices is then carefully normalized by dividing by `(len(data) - 1)`. The resulting value for each observation represents the exact cumulative proportion of the entire dataset observed up to that specific data point. This calculation is what transforms the sorted data points into cumulative probabilities.

The deliberate normalization process is crucial because it guarantees that the resulting CDF plot is correctly scaled, commencing at a cumulative probability near 0 for the minimum observation and culminating precisely at 1.0 for the maximum observation in the sample. The final output is a step function--although it often appears as a smooth, continuous curve when dealing with very large datasets--that visually maps the actual data values (plotted on the x-axis) directly to their corresponding cumulative probabilities (displayed on the y-axis). This visualization provides an

immediate, clear summary of the data spread and concentration.

Example 1: Generating the Empirical CDF for Raw Data Samples

Our first practical example demonstrates the core application of the eCDF methodology: calculating and plotting the cumulative distribution for a randomly generated dataset. This scenario mirrors a common real-world situation where a data professional receives raw, unstructured observations and requires an immediate, assumption-free visualization of its underlying distributional properties. For this implementation, we rely exclusively on the robust array manipulation capabilities of [NumPy](#) and the essential plotting features provided by [Matplotlib](#).

In this demonstration, we generate a large sample--specifically, 10,000 data points--drawn from a standard [Normal Distribution](#) using the function `np.random.randn(10000)`. Choosing a substantial sample size ensures that the resulting empirical CDF (eCDF) will closely mirror the known theoretical CDF of the normal distribution, thereby producing the familiar, characteristic S-shape curve expected from normal data. The comprehensive code block below illustrates the entire workflow, starting from the initial data generation through the final plotting of the cumulative curve:

```
import numpy as np
import matplotlib.pyplot as plt

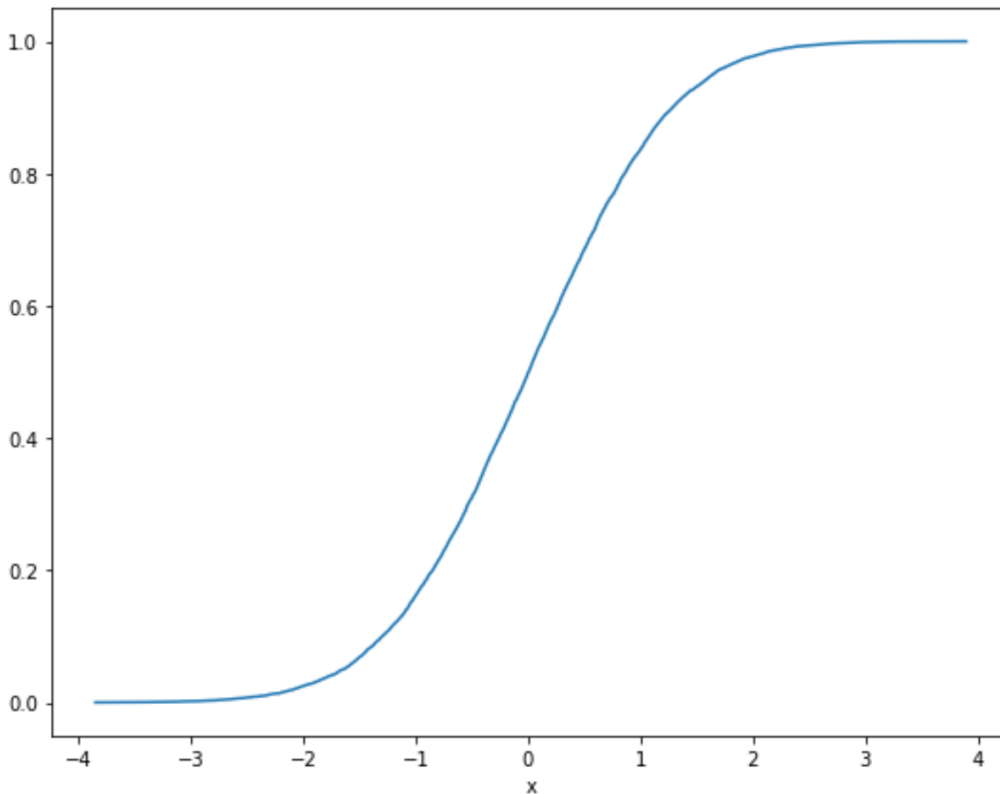
#define random sample of data
data = np.random.randn(10000)

#sort data
x = np.sort(data)

#calculate CDF values
y = 1. * np.arange(len(data)) / (len(data) - 1)

#plot CDF
plt.plot(x, y)
plt.xlabel('x')
```

The resulting output visualization effectively confirms the fundamental properties of the standard normal distribution. The **CDF** exhibits a slow, gradual increase at the extreme tails (indicating low data density) and then accelerates rapidly around the mean (signifying high data concentration). In the resulting plot shown below, the x-axis consistently represents the raw data values, while the y-axis precisely maps the corresponding cumulative probability values, ranging from 0 to 1.



Example 2: Leveraging SciPy for Theoretical Distribution CDFs

While the eCDF method is invaluable and highly adaptable for any empirical dataset, when the analyst possesses prior knowledge regarding the underlying distribution--such as knowing the data follows a [Normal Distribution](#), exponential distribution, or uniform distribution--it becomes significantly more efficient and mathematically accurate to calculate the theoretical [CDF](#). The powerful [SciPy](#) library, specifically its `scipy.stats` module, provides a comprehensive suite of built-in functions designed to calculate the CDF for virtually every standard theoretical distribution, bypassing the need for manual rank approximation.

To plot the theoretical CDF of a standard normal distribution, for example, we employ the dedicated function `scipy.stats.norm.cdf()`. This function accepts the sorted array of data points (our `x` array) as input and immediately returns the exact theoretical cumulative probability associated with those values, calculated strictly according to the distribution's defined parameters (mean and standard deviation). This specialized approach eliminates the need for the manual normalization and rank calculation inherent to the eCDF method, resulting in a perfectly smooth, mathematically precise curve.

The following code snippet demonstrates the integration of [SciPy](#) into the plotting workflow. It is important to note that while we still generate random data to establish the necessary range for the

x-axis (the domain of the plot), the corresponding y values are computed purely theoretically. This distinction ensures the resulting curve accurately represents the true population distribution function, rather than being merely an approximation based on the finite sample size.

```
import numpy as np
import scipy
import matplotlib.pyplot as plt

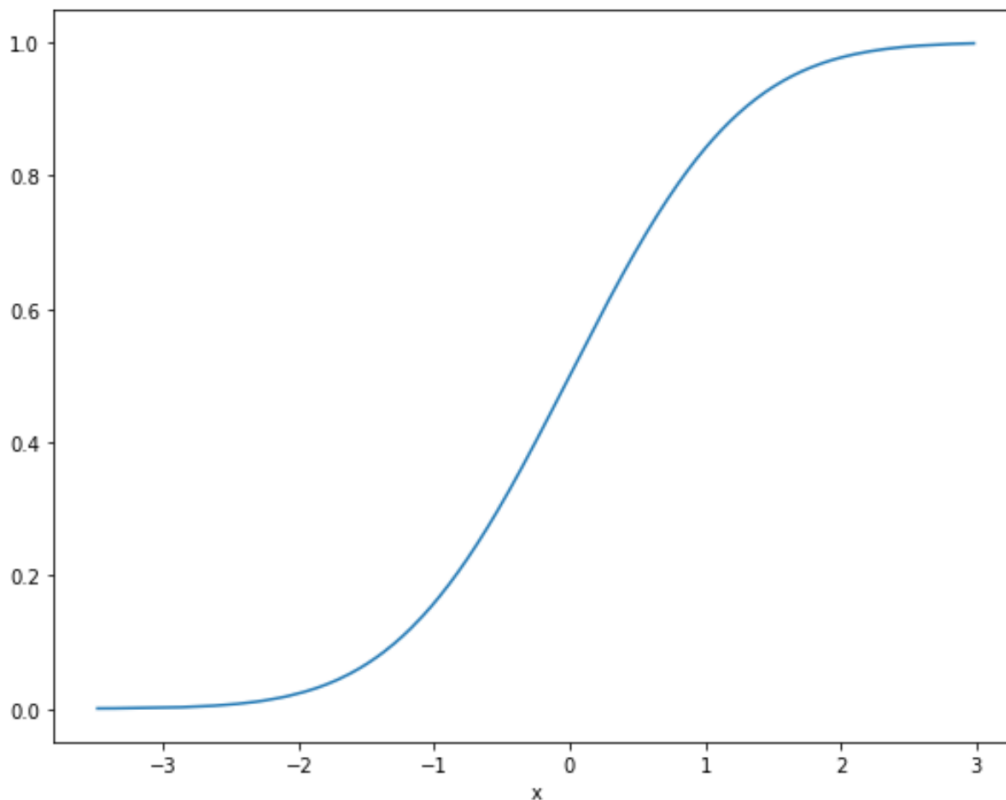
#generate data from normal distribution (to define x-range)
data = np.random.randn(1000)

#sort data
x = np.sort(data)

#calculate theoretical CDF values using scipy.stats
y = scipy.stats.norm.cdf(x)

#plot CDF
plt.plot(x, y)
plt.xlabel('x')
```

This theoretical modeling technique is strongly advocated whenever performing formal statistical inference where the population distribution is a critical assumption. It delivers a precise mathematical foundation rather than a sample-size-dependent approximation. The resulting theoretical plot, displayed below, appears visually analogous to the empirical plot but carries the rigorous accuracy derived directly from the mathematical properties of the distribution function.



Interpreting Quantiles and Distribution Comparison via the CDF

A calculated and plotted CDF offers substantial analytical advantages over alternative visualizations, such as standard histograms. Its most significant benefit lies in the straightforward, graphical identification of essential statistical metrics, specifically **percentiles and quantiles**. Because the y-axis strictly represents cumulative probability ranging from 0 to 1, analysts can immediately locate the specific data value (x-axis coordinate) below which a defined percentage of all observations fall. This direct mapping eliminates the need for complex calculations.

The process of extracting key descriptive statistics becomes remarkably simple. For instance, to identify the **median** of the dataset (the 50th percentile), an analyst merely locates the x-value corresponding to $y = 0.5$. Similarly, calculating the **interquartile range (IQR)** is achieved by determining the difference between the x-values corresponding to the third quartile ($y = 0.75$) and the first quartile ($y = 0.25$). This visual approach streamlines complex quantile analysis and provides rapid insights into the dataset's central tendency and spread.

Beyond single-dataset analysis, the CDF plot is an indispensable tool for comparing two or more distinct distributions simultaneously. By overlaying multiple CDF curves on the same axes, any observed vertical separation between the plots at a specific x-value immediately quantifies the difference in the cumulative probability (i.e., the difference in the proportion of data less than x).

between the two distributions. This principle of measuring separation between cumulative curves is the mathematical core of several powerful non-parametric statistical procedures, most notably the **Kolmogorov-Smirnov test**, which formally quantifies the maximum observed distance between two empirical or theoretical CDFs.

Advanced Resources for Statistical Distribution Analysis in Python

To solidify your expertise and maximize the application of statistical distributions within the [Python](#) environment, it is highly beneficial to engage with the extensive official documentation provided by the core scientific libraries. Achieving mastery over the high-performance array manipulation capabilities of [NumPy](#) and the diverse statistical toolset housed within [SciPy](#) is essential for transitioning from basic data plotting to sophisticated statistical modeling and inference workflows. These libraries offer optimization and flexibility far beyond the foundational examples covered here.

Key areas and tools that merit further exploration for comprehensive distribution analysis include:

Utilizing the fitting capabilities within `scipy.stats` to estimate and align distribution parameters (such as mean and variance) to match your observed empirical data.

Generating specialized **Probability Plots** (Q-Q plots) to visually and statistically assess how closely your data adheres to a specific theoretical distribution, such as the standard [Normal Distribution](#).

Applying advanced visualization techniques using [Matplotlib](#), often in conjunction with libraries like Seaborn, to customize CDF displays for publication-quality reports and enhanced analytical clarity.

Focusing on these advanced resources will facilitate a smooth progression from elementary data visualization techniques to full-scale statistical modeling and hypothesis testing.