

Learning the Empirical Rule: A Practical Guide with R

Authored by
Mohammed loot

November 2, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning the Empirical Rule: A Practical Guide with R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=8533>

The Foundation: Understanding the Empirical Rule in Statistical Analysis

The [Empirical Rule](#), often popularized as the 68-95-99.7 rule, stands as one of the most fundamental and intuitive principles in statistical theory. Its primary function is to offer a swift and highly reliable method for estimating the dispersion of data within a population, provided that the data rigorously adheres to a [normal distribution](#), commonly visualized as a bell curve. This powerful heuristic establishes a direct relationship between the percentage of data observations concentrated within specific intervals and the number of [standard deviations](#) those intervals are positioned away from the mean.

The ubiquity of the normal distribution across countless phenomena--ranging from human heights and test scores to measurement errors in physical sciences--is what lends the Empirical Rule its substantial predictive power. By acquiring just two parameters--the arithmetic mean (μ) and the standard deviation (σ)--we gain the ability to predict the precise concentration of observations. This capability is indispensable for core statistical tasks, including defining typical ranges, identifying statistical outliers, and establishing the practical significance of observed data points.

Formally, the Empirical Rule stipulates that for any dataset that perfectly follows a normal distribution:

68% of all data values are expected to fall within one standard deviation ($\mu \pm 1\sigma$) of the mean.

95% of all data values are expected to fall within two standard deviations ($\mu \pm 2\sigma$) of the mean.

99.7% of all data values are expected to fall within three standard deviations ($\mu \pm 3\sigma$) of the mean.

Throughout this comprehensive tutorial, we will transition from theory to practice by systematically applying the Empirical Rule using the statistical programming language **R**. We will leverage R's powerful built-in functions to numerically verify these theoretical percentages and demonstrate their application across varied, real-world datasets.

Leveraging the `pnorm()` Function for Distribution Calculations in R

To effectively analyze and manipulate the normal distribution within the R environment, statisticians rely heavily on a specific family of distribution functions. The function most central to implementing the Empirical Rule is `pnorm()`. This function is designed to calculate the value of the [cumulative density function](#) (CDF) for any specified value (quantile) within a normal distribution. In simple terms, `pnorm()` returns the cumulative area underneath the bell curve to the left of a given input point, representing the probability of observing a value less than or equal to that point.

Mastering the syntax of `pnorm()` is essential for unlocking its full utility in quantitative analysis. The function's basic structure is designed to be highly flexible, accommodating calculations based on the specific parameters of the distribution being analyzed:

pnorm(q, mean, sd)

The required arguments are defined as follows:

q: Represents the quantile, which is the specific numerical value of the normally distributed random variable for which we wish to evaluate the cumulative probability.

mean: Specifies the population arithmetic mean (μ) of the distribution. If this argument is omitted, R automatically defaults to a mean of 0, corresponding to the Standard Normal Distribution.

sd: Specifies the population standard deviation (σ) of the distribution. If this argument is omitted, R automatically defaults to a standard deviation of 1.

When the objective is to determine the probability that a value falls *between* two specific points (an upper bound and a lower bound), the method involves subtracting the CDF value of the lower bound from the CDF value of the upper bound. For applying the Empirical Rule, we are precisely interested in calculating this central area of probability concentrated symmetrically around the mean.

Numerical Verification of the 68-95-99.7 Rule Using R

We can use the `pnorm()` function to numerically confirm the theoretical underpinnings of the Empirical Rule. To simplify the verification process, we will assume the properties of the **Standard Normal Distribution**, where the mean (μ) is 0 and the standard deviation (σ) is 1. Under this standard assumption, one standard deviation corresponds to the values -1 and +1 on the x-axis.

To find the exact probability (or area) contained within one standard deviation of the mean, we calculate the cumulative area to the left of +1 and subtract the cumulative area to the left of -1. This subtraction technique effectively isolates and measures the central segment of the area under the curve.

The following R code snippet executes this calculation for one, two, and three standard deviations, confirming the highly precise percentages associated with the rule:

```
#find area under normal curve within 1 standard deviation of mean
```

```
pnorm(1) - pnorm(-1)
```

```
0.6826895
```

```
#find area under normal curve within 2 standard deviations of mean
```

```
pnorm(2) - pnorm(-2)
```

```
0.9544997
```

```
#find area under normal curve within 3 standard deviations of mean  
pnorm(3) - pnorm(-3)
```

```
0.9973002
```

A careful review of the output confirms the remarkable precision of the rule, showing values slightly more exact than the rounded integers typically quoted:

68.27% of data values fall within one [standard deviation](#) of the mean.

95.45% of data values fall within two standard deviations of the mean.

99.73% of data values fall within three standard deviations of the mean.

These calculations rigorously establish the theoretical basis of the rule within the R programming environment, laying the necessary groundwork for applying these principles to non-standard datasets that possess arbitrary means and standard deviations.

Example 1: Mapping the Empirical Rule to a Specific Dataset

In most realistic scenarios, data analysts rarely encounter the perfect standard normal distribution ($\mu=0$, $\sigma=1$). Instead, we must apply the Empirical Rule to datasets characterized by their own calculated means and standard deviations. This example illustrates the process of determining the specific numerical ranges that capture the 68%, 95%, and 99.7% thresholds for a custom distribution.

Imagine we are analyzing a dataset--perhaps measuring the lifespan of a product or the distribution of commute times--that is confirmed to be normally distributed with a mean (μ) of **7** and a [standard deviation](#) (σ) of **2.2**. Our goal is to identify the concrete numerical boundaries that define these key concentrations of data.

The required boundaries are determined by iteratively adding and subtracting multiples of the standard deviation (σ) from the mean (μ). For instance, the interval encompassing 68% of the data is mathematically defined by calculating $\mu \pm 1\sigma$. We execute these calculations directly in R, first clearly defining the mean and standard deviation variables for maximal code clarity and readability.

```
#define mean and standard deviation values
```

```
mean=7
```

```
sd=2.2
```

```
#find which values contain 68% of data (Mean +/- 1*SD)
```

```
mean-2.2; mean+2.2
```

```
4.8
```

9.2

#find which values contain 95% of data (Mean +/- 2*SD)

mean-2*2.2; mean+2*2.2

2.6

11.4

#find which values contain 99.7% of data (Mean +/- 3*SD)

mean-3*2.2; mean+3*2.2

0.4

13.6

The resulting output clearly delineates the exact ranges defined by the [Empirical Rule](#) specifically for this distribution:

68% of the data falls between **4.8** and **9.2** (the one standard deviation range).

95% of the data falls between **2.6** and **11.4** (the two standard deviation range).

99.7% of the data falls between **0.4** and **13.6** (the three standard deviation range).

These calculated boundaries are foundational for assessing the rarity of observations. For example, any value falling outside the 99.7% range--i.e., less than 0.4 or greater than 13.6--is statistically classified as an extremely rare event or an extreme outlier within the context of this specific normal distribution.

Example 2: Finding Percentage of Data Between Arbitrary Values

While the Empirical Rule is excellent for quick estimation based on 1, 2, or 3 standard deviations, researchers frequently need to calculate the precise percentage of data falling between two arbitrary values that do not align perfectly with these standard marks. This scenario demonstrates where the true flexibility and precision of the [pnorm\(\) function](#) are harnessed, enabling accurate probability calculations irrespective of the input values.

Let us analyze a [normal distribution](#), such as one representing IQ scores, known to have a mean (μ) of 100 and a standard deviation (σ) of 5. Suppose our interest lies in quantifying the proportion of the population whose scores are situated between **99** and **105**.

To determine this proportion, we must calculate the cumulative probability up to the upper bound (105) and then subtract the cumulative probability up to the lower bound (99). It is essential in this process to explicitly define the mean and standard deviation parameters within the `pnorm()` calls, ensuring the calculation references the correct distribution.

```
#find area under normal curve between 99 and 105
```

```
pnorm(105, mean=100, sd=5) - pnorm(99, mean=100, sd=5)
```

```
0.4206045
```

The resulting output, 0.4206045, is a direct measure of probability which translates immediately into a percentage. We can conclude with high precision that **42.06%** of the data observations fall between the values 99 and 105 for this specific IQ score distribution. This advanced application extends the power of the CDF far beyond the fixed thresholds of the standard Empirical Rule, enabling fine-grained, customized probability assessments.

Crucial Limitations and Context of the Empirical Rule

While the [Empirical Rule](#) is an intuitive and immensely useful statistical heuristic, it is imperative that analysts fully acknowledge its specific boundary conditions. The rule delivers strictly accurate results only when the underlying data distribution is demonstrably normal. Any significant deviations from true normality--such as pronounced skewness or the presence of heavy tails (leptokurtic distributions)--will inevitably lead to inaccurate estimations of data concentration.

Datasets that do not conform to the bell curve shape will not adhere to the 68-95-99.7 percentages. In such non-normal cases, the actual proportion of data falling within one or two [standard deviations](#) of the mean may differ substantially from the rule's predictions. Applying the Empirical Rule without confirming normality can therefore lead to misleading conclusions about the data spread and the identification of outliers.

For distributions whose shape is unknown or confirmed to be non-normal, statisticians must rely on a more conservative yet universally applicable theorem: **Chebyshev's Inequality**. This theorem provides a guaranteed minimum proportion of data that must fall within a certain number of standard deviations, irrespective of the distribution's shape. However, the bounds derived from Chebyshev's inequality are necessarily much wider than those provided by the Empirical Rule. This contrast underscores the significant efficiency and precision gained when the critical assumption of normality can be successfully confirmed through appropriate statistical tests (such as Q-Q plots or the Shapiro-Wilk test) prior to applying the Empirical Rule.

Further Resources for Advanced R Statistics

To solidify your expertise and further your understanding of distribution functions, probability theory, and statistical computing within the R environment, it is highly recommended to explore the following related topics and tools:

The full suite of R distribution functions, which includes `qnorm()` (the quantile function, for finding

values given a probability) and `dnorm()` (the density function, for plotting the curve).

Advanced methodologies for empirically testing the normality assumption of your collected data, such as the widely used Shapiro-Wilk test or visual inspection via Q-Q plots.

More complex applications of the [pnorm\(\) function](#) in calculating critical statistical measures like confidence intervals, statistical power, and p-values for hypothesis testing.

A firm mastery of these fundamental statistical concepts and their implementation in R provides an exceptionally strong analytical foundation required for advanced statistical modeling and sophisticated data analysis projects.