

# How to Understand and Work with Open Ended Distributions

Authored by  
**stats writer**

December 5, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Understand and Work with Open Ended Distributions*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=105783>

In the field of statistics, understanding how data is grouped and presented is fundamental to accurate analysis. An **open ended distribution** is a specific type of frequency distribution characterized by the presence of one or more classes (or "bins") that lack a defined upper or lower boundary. This open-ended nature means that while we know the frequency of observations falling into these extreme categories, we lose precision regarding the exact value of those observations. Such distributions are common when dealing with datasets that exhibit extremely low or extremely high outlier values, or when survey design intentionally limits the specificity of responses at the tails of the data range.

The presence of open ends significantly impacts the methods we can use to summarize and interpret the dataset. For instance, if the lowest class is defined as "Under \$10,000," we know how many individuals fall into that category, but we do not know if their income is \$1,000 or \$9,999. Similarly, if the highest class is labeled "Over \$500,000," the actual income could be slightly above that threshold or substantially higher, creating ambiguity. This lack of precise boundary information distinguishes open ended distributions from their closed-ended counterparts and necessitates the use of robust, non-parametric analytical techniques, particularly focusing on measures like the median rather than the mean.

## Defining Open vs. Closed Distributions

To fully appreciate the implications of an open ended distribution, it is helpful to contrast it with a standard, **closed ended distribution**. In a closed ended distribution, every single class interval possesses both a clearly defined upper limit and a clearly defined lower limit. For example, income classes might be delineated as "\$20,000 - \$39,999," "\$40,000 - \$59,999," and so on, with no ambiguity regarding the endpoints of the range. This structure allows researchers to confidently assign a midpoint (or class mark) to every interval, which is crucial for calculating descriptive statistics like the estimated arithmetic mean and the standard deviation.

Conversely, an open ended distribution features at least one interval where one boundary is missing. This usually occurs at the distribution's extremities. If the class is the smallest, it is typically left open at the lower end (e.g., "Less than 10 years"). If the class is the largest, it is left open at the upper end (e.g., "70 years and older"). It is possible, though less common, for a distribution to be open at both ends simultaneously. This structural characteristic, while often necessary for practical data collection, fundamentally alters the mathematical tractability of the resulting data presentation.

Consider the following examples illustrating how open boundaries appear in tabulation. The first image demonstrates a scenario where the smallest class interval is open ended:

Annual Income	Frequency
<\$20,000	6
\$20,000 < \$39,999	12
\$40,000 < \$59,999	17
\$60,000 < \$79,999	19
\$80,000 < \$99,999	14
\$100,000 < \$119,999	4

The next example showcases an open ended distribution where the largest class interval lacks a defined upper limit, often used for variables like wealth or income where values can increase indefinitely:

Annual Income	Frequency
\$10,000 < \$20,000	6
\$20,000 < \$39,999	12
\$40,000 < \$59,999	17
\$60,000 < \$79,999	19
\$80,000 < \$99,999	14
> \$100,000	4

For comparison, a classic closed ended distribution ensures every respondent is classified within explicit boundaries:

Annual Income	Frequency
\$10,000 < \$20,000	6
\$20,000 < \$39,999	12
\$40,000 < \$59,999	17
\$60,000 < \$79,999	19
\$80,000 < \$99,999	14
\$100,000 < \$119,999	4

## Why Researchers Employ Open Ended Classes

The decision to implement an open ended class structure is typically driven by practical considerations in data collection, rather than mathematical necessity. Open ended distributions are frequently the result of researchers deliberately designing surveys in a way that minimizes participant discomfort and maximizes the response rate. When dealing with sensitive topics, such as annual household income, extreme age, or net worth, individuals at the very high or very low ends of the scale may be reluctant to disclose precise figures.

For instance, imagine a large-scale demographic survey. A researcher might choose to define the maximum income response as "Over \$200,000." This approach serves two key purposes. Firstly, it provides a degree of privacy, allowing high-earning respondents to participate without revealing proprietary or highly specific financial data. Secondly, it manages data heterogeneity; often, the precise differences between extremely high values (e.g., \$500,000 vs. \$1,000,000) may not be statistically relevant to the overall research question, especially if the focus is on the middle 90% of the population.

Similarly, open-ended classes can efficiently handle outliers. If a small number of observations fall far outside the typical range, creating a separate, precise class for each outlier would result in a scattered and difficult-to-read frequency distribution. Grouping these extreme values into a single open-ended class helps maintain the clarity and interpretability of the distribution table while acknowledging the presence of exceptional data points. Therefore, in essence, researchers utilize open ended classes as a strategic tool to enhance the practicality and ecological validity of their data collection processes.

## The Statistical Limitations of Open Ended Data

While open ended distributions offer practical benefits in survey administration, they introduce significant mathematical challenges, primarily due to the phenomenon known as censored data. When data is censored, we lose access to the true, underlying numerical values within the open interval. We know the count (frequency) within the class, but we do not know the magnitude of the individual data points. This ambiguity makes it impossible to calculate exact descriptive statistics that rely on every data point's precise value.

The primary casualty of open-ended classification is the accurate calculation of the mean (arithmetic average) and the standard deviation. To calculate the mean from a grouped frequency distribution, one typically estimates the midpoint of each class and uses that midpoint to represent all observations within that class. However, for an open class like "> \$100,000," there is no mathematically justifiable midpoint. Should we assume the average income is \$150,000, \$200,000, or perhaps \$500,000? Any arbitrary assignment of a midpoint introduces bias and significantly compromises the validity of the calculated mean.

Furthermore, the inability to calculate an accurate mean also prevents the calculation of the exact standard deviation, which measures data dispersion around the mean. Since the standard deviation requires precise measurement of the deviation of each class midpoint from the mean, the lack of a reliable class mark for the open interval renders the final measure of variability highly unreliable. Consequently, researchers must turn to alternative, more robust measures of central tendency and variability that are less sensitive to extreme, unknown values, thereby preserving the integrity of their statistical summary.

## Alternative Measures of Central Tendency

Since the mean is rendered unreliable in the presence of open ended classes, statisticians rely heavily on the **median** and the **mode** as primary measures of central tendency. The median, defined as the middle value of the dataset when ordered, is particularly useful because its calculation depends only on the cumulative frequencies and the location of the median group, not on the specific magnitudes of the values in the extreme open-ended classes.

The mode, which represents the most frequently occurring value or class, is also a viable measure. It is determined simply by observing the class with the highest frequency, regardless of whether that class is open or closed. However, the mode only provides information about the most common outcome, which may not accurately reflect the overall center of the distribution if the data is highly skewed. For most analytical purposes, the estimated median offers a far more representative summary of the typical value in an open ended distribution.

Crucially, the median's robustness stems from its positional nature. As long as the open ended class does not contain the median observation itself--meaning the 50th percentile falls within a closed interval--the existence of the open ends does not affect the calculation. Even if the median falls within a closed class, we can accurately estimate its value using a specific formula designed for grouped data. This makes the median the standard, preferred statistic when analyzing data presented in open ended frequency distributions.

## Calculating the Estimated Median for Grouped Data

To obtain the best possible summary of the center of an open ended distribution, we must calculate the **estimated median**. This process involves two main steps: first, identifying the median class (the interval containing the  $N/2$  observation), and second, interpolating the exact median value within that class using the cumulative frequencies. This interpolation method assumes that the values within the median class are evenly distributed, allowing for a precise estimation despite the data being grouped.

The formula used to find the best estimate of the median is standard practice in non-parametric statistics when faced with grouped data:

**Best Estimate of Median:**  $L + (n/2 - F) / f * w$

Each variable in this formula represents a specific characteristic derived from the frequency table, ensuring that the estimated median accurately reflects the cumulative progression of the data:

**L:** Represents the **lower limit** of the median group (the class interval where the median falls). This boundary is essential as it forms the starting point for interpolation.

**n:** Denotes the **total number of observations** in the entire dataset, which determines the exact position of the middle value ( $n/2$ ).

**F:** Refers to the **cumulative frequency** of all classes preceding the median group. This count tells us how many observations occur before the median class begins.

**f:** Is the **frequency** of the median group itself, indicating the density of observations within that critical interval.

**w:** Represents the **width** of the median group, calculated as the difference between its upper and lower boundaries.

By utilizing this formula, we effectively place the estimated median at the proportional point within the median class that corresponds to the 50th percentile of the total observations.

## A Practical Example of Median Estimation

To illustrate the calculation, let's use the earlier example of an open ended income distribution where the highest class is open:

Annual Income	Frequency
\$10,000 < \$20,000	6
\$20,000 < \$39,999	12
\$40,000 < \$59,999	17
\$60,000 < \$79,999	19
\$80,000 < \$99,999	14
> \$100,000	4

First, we determine the total number of observations, which in this case, is 72. The median observation position is therefore  $N/2$ , or  $72/2 = 36$ . We need to find the class interval that contains both the 36th and 37th observations (for an even dataset). By calculating the cumulative frequencies (e.g., if the first class has 10, the second 15, the cumulative frequency before the third class is 25), we can pinpoint the median class.

In this scenario, the median class is identified as "\$60,000 - \$79,999." Before this class, the

cumulative frequency (F) is 25. The frequency (f) of this specific class is 19. The lower limit (L) is \$60,000, and the class width (w) is \$79,999 - \$60,000 + \$1 (or simply the class width if dealing with continuous data), which is \$19,999.

Plugging these values into the interpolation formula yields the estimated median:

$$\text{Median: } 60,000 + ( (72/2 - 25) / 19 ) * 19,999$$

$$\text{Median: } 60,000 + ( (36 - 25) / 19 ) * 19,999$$

$$\text{Median: } 60,000 + (11 / 19) * 19,999 \approx 60,000 + 11,578$$

The resulting best estimate of the median annual income for individuals in this dataset is **\$71,578**. This robust figure provides a reliable measure of the center, unaffected by the unknown precise values in the open-ended class "≥ \$100,000."

## Mitigating the Impact of Open Ended Data

While the use of the median successfully bypasses the need for precise values in open classes, researchers often seek methods to estimate or bound the true mean and standard deviation for more comprehensive analysis. One common technique is to use a statistical imputation method, where a reasonable value is assigned to the open class based on external data or theoretical distribution assumptions. For instance, if external studies suggest incomes over \$100,000 typically follow a Pareto distribution, researchers might use parameters from that distribution to estimate the mean income within the open interval.

Another approach, particularly useful when the open class is at the upper end of a monetary or demographic variable, is to set an educated, conservative upper bound. For example, in an income study, while the category is "Over \$200,000," the researcher might tentatively cap the class midpoint calculation at \$350,000, acknowledging that this is an estimation, but one that allows for the calculation of an \*estimated\* mean and standard deviation. It is crucial, however, that any such assumption be clearly documented and justified, as it fundamentally involves making assumptions about the censored data.

Finally, researchers can explore non-parametric statistical tests, which do not rely on assumptions of normality or precise measures of the mean and standard deviation. Tests based on ranks (like the Wilcoxon rank-sum test or the Kruskal-Wallis test) utilize the ordinal nature of the data, making them appropriate for comparisons involving open ended distributions. By carefully selecting analytical tools and transparently documenting any necessary assumptions, the limitations imposed by open ended data can be managed effectively, ensuring that the findings remain scientifically sound.