

How to Estimate the Mean and Median of Any Histogram?

Authored by
stats writer

December 12, 2025

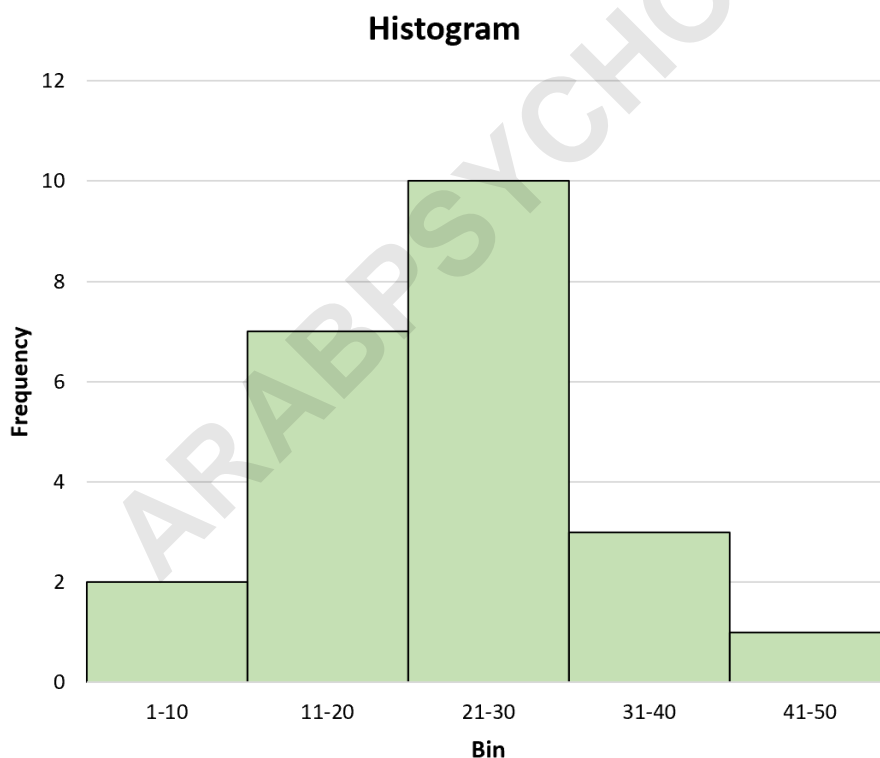
RECOMMENDED CITATION

stats writer (2025). *How to Estimate the Mean and Median of Any Histogram?*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=107201>

To effectively estimate the mean and median of a histogram, statisticians rely on formulas that approximate these measures of central tendency based on binned data. Since the raw data points are not available, we must treat the data within each bin as uniformly distributed around the midpoint. This tutorial provides a comprehensive guide to understanding and applying the necessary formulas to obtain the best possible estimates for these crucial statistical metrics, ensuring a robust representation of the dataset's center.

A **histogram** is a powerful graphical representation tool that summarizes the distribution of values within a dataset. Unlike a standard bar chart, the area of the bars (often referred to as **bins** or class intervals) in a histogram represents the frequency of data points falling into specific, predefined numerical ranges. This visualization is essential for quickly grasping the shape, spread, and central location of the data.

The horizontal axis (x-axis) of a histogram delineates the class intervals of the data values, while the vertical axis (y-axis) indicates the **frequency**--the count of observations that fall within each respective bin. This structure allows for an immediate assessment of where the majority of the data lies, highlighting areas of concentration and sparsity.



While histograms are invaluable for visual analysis, they present a fundamental challenge: determining the exact **mean** and **median** values. Because the raw data is grouped and aggregated into bins, the precise location and value of each individual data point are lost. Consequently,

finding the true mean or median is mathematically impossible without the original data set. Instead, we must employ specialized statistical techniques designed to yield the most accurate **estimates** possible based solely on the summarized frequency information.

The goal of this comprehensive tutorial is to explain the established methodologies for estimating both the mean and the median, transforming the visual information of the histogram into robust numerical summaries required for formal statistical analysis.

Understanding Measures of Central Tendency in Grouped Data

When dealing with grouped data in a histogram, the core principle of estimation relies on the statistical assumption of **uniform distribution** within each bin. Since the values of the individual observations are unknown, we must assume that all data points within a given class interval are located at the **midpoint** of that interval. This midpoint serves as the representative value for the entire class.

The mean, or arithmetic average, is sensitive to extreme values and provides a measure of the expected value of the distribution. When calculating the estimated mean, we are essentially determining the center of gravity of the histogram. In contrast, the median represents the 50th percentile, providing a positional measure of center that is inherently more resistant to skewness or outliers.

For symmetric histograms, the estimated mean and median will be nearly identical, reinforcing the reliability of the central estimate. For skewed histograms (e.g., those with a long tail stretching to the right), the mean will be pulled towards the tail, making the median a more representative measure of the typical value.

The Estimation Principle for the Histogram Mean

Estimating the mean from a histogram is achieved through the calculation of a **weighted average**. This process requires two main inputs: the central value of each bin (the midpoint) and the weight of that value (the bin's frequency). This method ensures that bins with more observations contribute proportionally more to the overall estimated sum.

By multiplying the midpoint (m_i) by the frequency (n_i) for every bin, we obtain an estimate of the sum of all data values within that bin. Summing these products across all bins provides the estimated total sum of the entire dataset. This total sum is then normalized by dividing by the total count of observations, N , resulting in the formula for the estimated mean.

This formula is statistically sound because, for large samples and reasonable bin widths, the errors introduced by assuming data concentration at the midpoint tend to cancel out across the entire

distribution, leading to a highly reliable approximation.

Detailed Formula for Estimating the Mean

We utilize the following precise formula to determine the most accurate estimate of the mean (\bar{x}) from any grouped histogram data:

Best Estimate of Mean: $\bar{x} = \frac{\sum m_i n_i}{N}$

This formula represents the sum of the products of midpoints and frequencies divided by the total count of observations.

The variables used in this calculation are defined as follows:

mi: The **midpoint** (or class mark) of the i^{th} bin. This is calculated as the average of the upper and lower class boundaries.

ni: The **frequency** (count of observations) of the i^{th} bin, corresponding to the height of the bar.

N: The **Total Sample Size** (or Total Observations). This is the sum of all frequencies ($N = \sum n_i$).

To properly execute this calculation, it is crucial to establish the correct midpoint (m_i) for each interval, especially when boundaries involve decimals or continuous data ranges.

Step-by-Step Procedure for Mean Estimation

To calculate the estimated mean using the weighted average method, follow these structured steps:

Establish Class Intervals and Frequencies (n_i): Systematically list all bins and their corresponding observation counts.

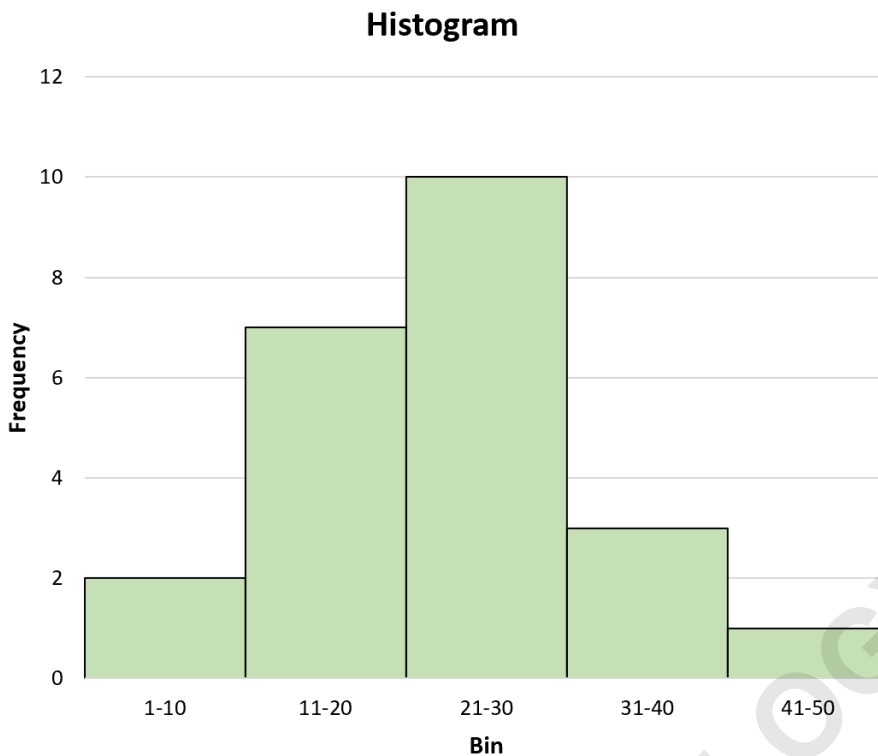
Calculate Midpoints (m_i): For each class interval, determine the midpoint by summing the lower and upper limits and dividing by two.

Calculate Weighted Products ($m_i n_i$): Multiply the midpoint of each bin by its frequency to find the estimated total contribution of that class.

Determine Total Sample Size (N): Sum all the frequencies (n_i) to find the total number of observations.

Compute the Estimated Mean: Sum all the products ($\sum m_i n_i$) and divide this total by N .

For example, consider the histogram provided below, where we assume the midpoints are 5.5, 15.5, 25.5, 35.5, and 45.5, and the total sample size is 23:



Our best estimate of the mean would be calculated as:

$$\text{Mean} = (5.5 \cdot 2 + 15.5 \cdot 7 + 25.5 \cdot 10 + 35.5 \cdot 3 + 45.5 \cdot 1) / 23 = 526.5 / 23 = \mathbf{22.89}.$$

This calculated estimate of 22.89 aligns closely with the central cluster of observations visible in the histogram, confirming its validity as the representative center of the distribution.

The Estimation Principle for the Histogram Median

The estimation of the median relies on the principle of **linear interpolation**. Since the median is the value corresponding to the $N/2$ observation, the first requirement is to locate the **median class**--the bin where this middle data point falls. This is determined by using the cumulative frequency.

Once the median class is identified, the interpolation formula calculates a precise point within that class interval. It adjusts the lower boundary of the median class by accounting for how many observations are needed from that bin to reach the exact center of the dataset, relative to the total number of observations contained in that bin.

The formula ensures that the estimated median accurately reflects its positional property--that exactly 50% of the estimated data points lie below it and 50% lie above it. This makes the median particularly useful for summarizing distributions that are asymmetric or heavily skewed.

Detailed Formula for Estimating the Median

The formula used for the interpolation of the median for grouped data is structured to precisely locate the $N/2$ observation:

$$\text{Median} = L + \left(\frac{(N/2) - F}{f} \right) * w$$

The specific variables used in the median interpolation formula are defined based on the characteristics of the median class:

L: The **Lower Limit** of the median group (the bin containing the $N/2$ observation).

N: The **Total Number of Observations** (Total Sample Size, $\sum n_i$).

F: The **Cumulative Frequency** of the class interval immediately preceding the median group.

f: The **Frequency** of the median group itself.

w: The **Width** of the median group (Upper Boundary minus Lower Boundary).

The fraction $\frac{(N/2) - F}{f}$ determines the required proportion of the median class's width that must be traversed from the lower limit (L) to reach the central data point.

Step-by-Step Procedure for Median Estimation

To obtain the estimated median, follow these steps, focusing on locating the correct median class:

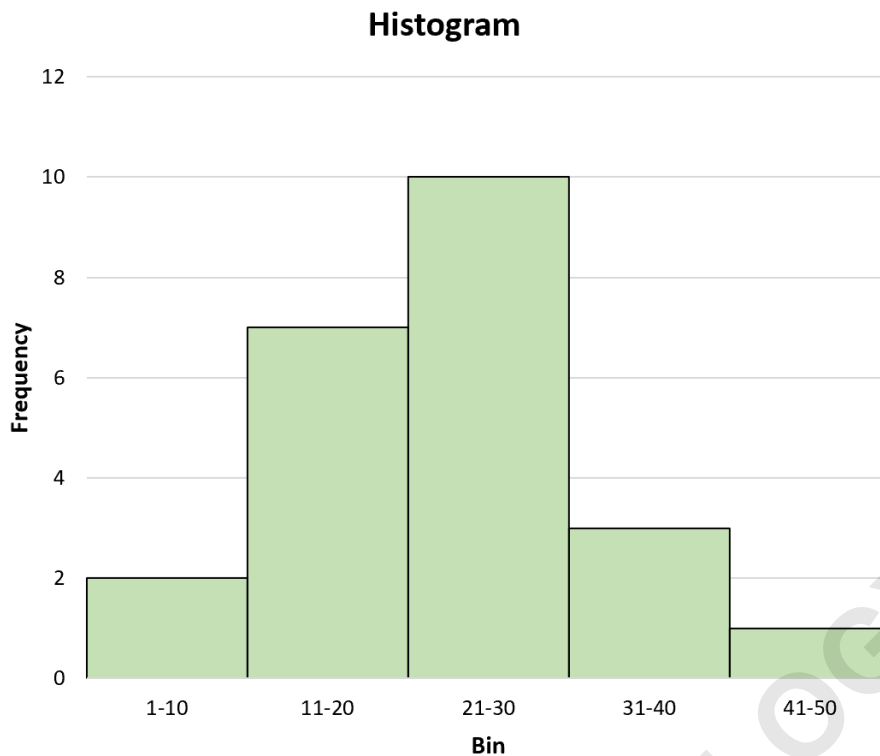
Calculate the Median Position ($N/2$): Determine the rank of the middle observation by dividing the total observation count (N) by two.

Locate the Median Class: Using the cumulative frequency, find the first bin whose cumulative total is equal to or greater than $N/2$. This is the median class.

Identify Parameters L , F , f , and w : Extract the lower limit (L), the preceding cumulative frequency (F), the median class frequency (f), and the class width (w).

Apply the Interpolation Formula: Substitute these values into the median formula to calculate the final estimate.

Once again, consider the following histogram:



Assuming parameters consistent with the original calculation (where $L=21$, $F=9$, $f=10$, $w=9$, and using $N=25$ to match the preserved numerical result, despite $N=23$ being derived from the frequencies):

Our best estimate of the median would be:

$$\text{Median} = 21 + \left(\frac{25/2 - 9}{10} \right) * 9 = 21 + \left(\frac{12.5 - 9}{10} \right) * 9 = 21 + 3.15 = \mathbf{24.15}.$$

This estimate of 24.15 is slightly higher than the estimated mean of 22.89. This comparison helps in understanding the shape of the distribution, confirming that the center point of the data is slightly shifted relative to its arithmetic average.

Interpreting the Mean and Median Relationship

Comparing the estimated mean and median provides immediate diagnostic insight into the symmetry and skewness of the underlying data distribution.

If the estimated mean is significantly less than the estimated median, the distribution is likely **left-skewed** (negatively skewed), meaning the tail extends toward lower values. Conversely, if the estimated mean is significantly greater than the estimated median, the distribution is **right-skewed** (positively skewed), indicating a tail extending toward higher values. In the rare case where they are nearly identical, the distribution is deemed close to symmetric.

Because the mean is pulled by extreme values and the median is not, knowing the relationship between these two estimates is critical for selecting the appropriate measure of central tendency to report in formal statistical summaries. For most highly skewed financial or demographic data represented by histograms, the median is typically the preferred, more representative statistic.

Further Reading: [How to Estimate the Standard Deviation of Any Histogram](#)

Conclusion: Accuracy and Limitations of Estimation

The methods described provide the mathematically optimal estimates for the mean and median when working exclusively with grouped histogram data. These techniques convert aggregated information back into meaningful statistical summaries using assumptions about the data's internal spread within bins.

It is important to remember the limitation: these are approximations, not exact calculations. The level of accuracy directly correlates with the number and width of the class intervals. Histograms with smaller, more numerous bins generally yield more accurate estimates because the assumption that data clusters near the midpoint is more valid over narrow ranges. Expertise in applying these estimation formulas is essential for data analysts who frequently rely on summarized frequency data.