

# Learning to Control Histogram Bin Sizes Using SAS

Authored by  
**Mohammed looti**

May 12, 2026

## RECOMMENDED CITATION

Mohammed looti (2026). *Learning to Control Histogram Bin Sizes Using SAS*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=3596>

## Controlling Data Visualization: Specifying Bins in SAS Histograms

When conducting [data visualization](#), [histograms](#) are vital instruments used to understand the frequency distribution of numerical [variables](#). A key factor in producing an insightful histogram is the accurate definition of its [bins](#)--the continuous intervals that group the raw data points. Within the powerful statistical software [SAS](#), the specialized ``midpoints`` statement, used inside the ``PROC UNIVARIATE`` procedure, grants the analyst granular control over these bin definitions. This control allows for precise tailoring of the visualization to achieve optimal analytical clarity and insight. This comprehensive guide will detail the essential steps required to explicitly define both the count and the width of bins when generating histograms in SAS.

Relying solely on the automated binning methods chosen by default in [SAS](#) can sometimes mask important underlying patterns or subtle anomalies present within your data distribution. This is because default settings prioritize efficiency and generality, which may not always align with specific research questions. By proactively setting the ``midpoints`` for your [histogram](#), you gain the ability to fine-tune its visual granularity. This process ensures the resulting visualization is highly informative, accurate, and perfectly aligned with your specific analytical objectives. This high degree of flexibility is indispensable, whether you are conducting preliminary [exploratory data analysis](#) or preparing findings for critical communication to stakeholders.

## Understanding Histograms and the Significance of Bins

Fundamentally, a [histogram](#) serves as a robust graphical representation illustrating the probability distribution of a continuous numeric [variable](#). It is constructed using a series of adjacent rectangular bars. Each bar, critically known as a [bin](#), corresponds to a specific, non-overlapping range of data values. The height of the bin is directly proportional to the frequency or percentage of [observations](#) that fall precisely within that designated range. Conventionally, the horizontal [x-axis](#) of the chart displays the range of data values (the variable itself), while the vertical [y-axis](#) quantifies the count or proportion of data points.

The methodology employed to group the data into [bins](#) exerts a profound influence on the ultimate visual shape and subsequent interpretability of the [histogram](#). If the bins are set to be overly wide (resulting in an insufficient number of bins), crucial structural details about the data's inherent shape--such as signs of bimodality, significant skewness, or the emergence of distinct data clusters--can be tragically obscured or completely averaged out. Conversely, employing bins that are excessively narrow (leading to an overwhelmingly large number of bins) can cause the histogram to appear overly noisy and erratic. This noise makes it exceptionally challenging to distinguish legitimate overall trends from mere random sampling fluctuations or outliers.

Therefore, the deliberate and thoughtful selection of the optimal number and width of [bins](#)

constitutes a foundational requirement for generating truly insightful [data visualization](#). This mastery allows analysts to thoroughly explore various facets of the data's distribution, accurately identify central tendencies, reliably assess the spread or variance of the data, and confidently detect potential outliers. In the context of [SAS](#) programming, the `midpoints` statement furnishes users with direct, explicit control over this critical grouping process, enabling them to bypass potentially misleading default settings and instead generate visualizations that precisely support their complex analytical objectives.

## Leveraging `PROC UNIVARIATE` with the `MIDPOINTS` Statement

Within the [SAS](#) environment, the `PROC UNIVARIATE` procedure stands out as an exceptionally versatile command. It is meticulously engineered for calculating descriptive statistics, performing tests of location, and producing a comprehensive array of graphical outputs, including highly customizable [histograms](#). Although `PROC UNIVARIATE` possesses sophisticated internal algorithms designed to automatically suggest an appropriate binning strategy based on common statistical rules, the power of the `midpoints` option lies in its capacity to offer the user complete, explicit control over this binning process. This option is required to be specified directly within the `histogram` statement, where it precisely defines the center point for every single bin to be displayed.

The standard syntax structure required for effectively employing the `midpoints` statement to customize the bin arrangement of your histogram is illustrated in the code block below. Understanding this structure is fundamental to applying custom binning strategies to any [variable](#) within your [dataset](#).

```
proc univariate data=my_data;  
histogram my_variable / midpoints=(9 to 36 by 3);  
run;
```

In this foundational SAS syntax, `my_data` serves as the reference to your specific input [dataset](#), and `my_variable` strictly denotes the continuous numeric [variable](#) whose distribution you intend to visualize. The `midpoints` option requires a specific range definition, structured as `(start to end by interval)`. Here, `start` dictates the central value of the very first bin, `end` establishes the central value of the final bin, and most importantly, `interval` determines the fixed, consistent distance between the center points of adjacent bins. This `interval` parameter holds the direct key to controlling the physical width of each [bin](#), consequently defining the total count of bins presented in the resulting histogram. For instance, the specification `(9 to 36 by 3)` will construct bins centered at 9, 12, 15, and so forth, up to 36, ensuring that every bin maintains a uniform width of 3 units.

It is absolutely imperative to strategically select a ``start`` value that is less than or equal to the minimum value observed within your [dataset](#), and an ``end`` value that is greater than or equal to the dataset's maximum value. This careful boundary selection guarantees that all [observations](#) are correctly encompassed and represented within the final histogram visualization. Furthermore, the ``interval`` parameter fundamentally governs the level of granularity: a smaller interval generates a greater quantity of narrower bins, providing a highly detailed visualization; conversely, a larger interval yields fewer, wider bins, which serves to offer a broader, more generalized overview of the data distribution.

## Practical Example: Creating and Analyzing a Default Histogram

To properly illustrate the practical implications of utilizing the ``midpoints`` statement, we must first establish a robust sample [dataset](#) in the [SAS](#) environment. For this demonstration, we will create a small dataset containing realistic, illustrative information concerning basketball players, specifically tracking metrics such as their team affiliation, points scored, and rebounds collected. This intentionally structured data will serve as the controlled foundation for all subsequent histogram examples, enabling us to clearly compare and contrast how diverse binning strategies dramatically alter the final visualization and interpretation.

The following [SAS](#) program block meticulously creates a dataset named ``my_data`` and populates it with a defined set of fourteen player [observations](#) using the ``DATALINES`` statement. Following the data creation step, we utilize ``PROC PRINT`` to display the entire contents of this dataset. This step acts as a crucial verification checkpoint, confirming that the data structure is correct and that the dataset is fully prepared for subsequent statistical analysis and [visualization](#) routines.

```
/*create dataset*/  
data my_data;  
input team $ points rebounds;  
datalines;  
A 29 8  
A 23 6  
A 20 6  
A 21 9  
A 33 14  
A 35 11  
A 31 10  
B 21 9  
B 14 5  
B 15 7  
B 11 10
```

```
B 12 6
B 10 8
B 15 10
;
run;

/*view dataset*/
proc print data=my_data;
```

Upon successful execution of the code above, [SAS](#) will output a formatted table that clearly showcases the contents of our newly created ``my_data`` dataset, confirming the successful load of all player statistics. This visual confirmation step is essential for proceeding with confidence to the analytical phase.

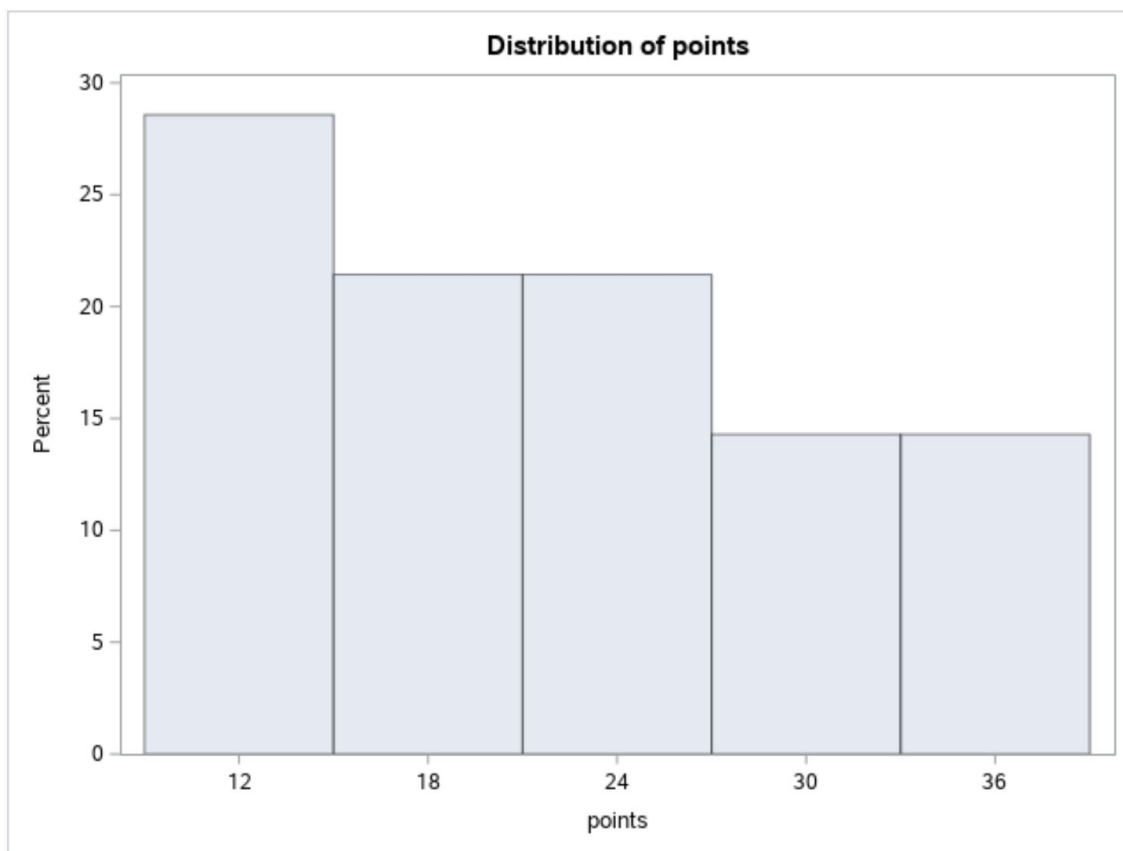
Obs	team	points	rebounds
1	A	29	8
2	A	23	6
3	A	20	6
4	A	21	9
5	A	33	14
6	A	35	11
7	A	31	10
8	B	21	9
9	B	14	5
10	B	15	7
11	B	11	10
12	B	12	6
13	B	10	8
14	B	15	10

Our next critical step is to generate a preliminary, standard [histogram](#) for the ``points`` [variable](#), strictly relying on the default binning algorithms inherent in ``PROC UNIVARIATE``. This initial visualization will serve as a definitive baseline against which we can compare the subsequent customized histograms, effectively demonstrating how SAS automatically calculates and implements bin widths when the user provides no explicit ``midpoints`` instruction.

The following concise [SAS](#) code block is used to produce the default histogram for the ``points`` variable sourced from our ``my_data`` dataset:

```
/*create histogram for points variable*/  
proc univariate data=my_data;  
histogram points;  
run;
```

Upon execution, [SAS](#) generates a histogram similar to the image presented below. The horizontal [x-axis](#) delineates the full range of `points` scored by the players, while the vertical [y-axis](#) accurately represents the calculated percentage of [observations](#) that are contained within each automatically defined bin.



A careful examination of this default [histogram](#) reveals that the SAS software has automatically chosen bin midpoints that are consistently separated by intervals of 6 units. While this provides a rudimentary, general overview of the data's central tendency and spread, it often falls short of being the most analytically insightful representation for all research objectives. Key fine details of the distribution may be inadvertently smoothed over, or, alternatively, distinct clusters of data points might be combined into a single, less informative, broader bin. This observation emphatically highlights the critical analytical value of exerting explicit, manual control over the histogram's bin width.

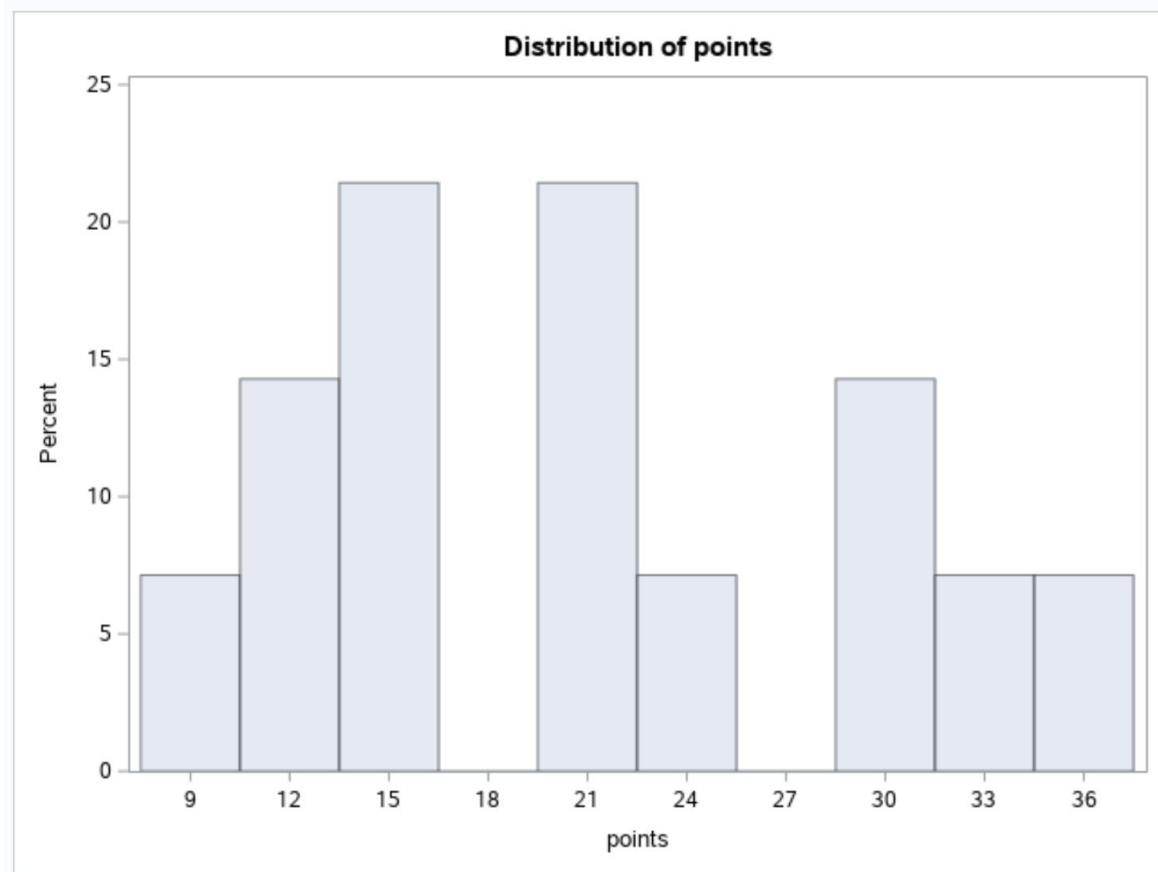
## Customizing Bins for Different Analytical Needs

The strategic ability to customize the number and width of histogram [bins](#) is paramount for exploring a dataset at various, targeted levels of analytical detail. Depending entirely on the scope of your analytical goals, you may strategically decide to increase the number of bins to fully expose finer, subtle data structures, or conversely, decrease the number of bins to achieve a more generalized, high-level overview. The ``midpoints`` statement within the ``PROC UNIVARIATE`` procedure provides precisely this necessary flexibility and control over the visualization output.

To achieve a more granular and detailed understanding of the distribution characteristics of the ``points``, we must effectively increase the total number of [bins](#). This is accomplished by specifying a significantly smaller ``interval`` value for the ``midpoints`` option. A reduced interval inherently results in narrower individual bins, which consequently means more bins fit across the same overall data range. This process successfully reveals the finer, often hidden structures within the data. Let us modify our previous [SAS](#) code to explicitly define ``midpoints`` that occur at fixed intervals of 3, spanning the entire observed range from 9 to 36. This calculated choice will dramatically increase the level of detail compared to the SAS default setting analyzed previously.

```
/*create histogram for points variable with custom bins (Interval=3)*/  
proc univariate data=my_data;  
  histogram points / midpoints=(9 to 36 by 3);  
run;
```

The resulting [histogram](#), presented immediately below, clearly exhibits a much greater number of bars than the default version. Each individual [bin](#) now encompasses a significantly narrower range of ``points``, which permits the analyst to observe more subtle peaks, valleys, and potential micro-clusters in the data distribution. This enhanced level of detail is particularly invaluable when the objective is to precisely identify specific score groupings or gaps in performance that might have been entirely masked or smoothed over when using the much broader default bins.



## Exploring Generalized Overviews with Wider Bins

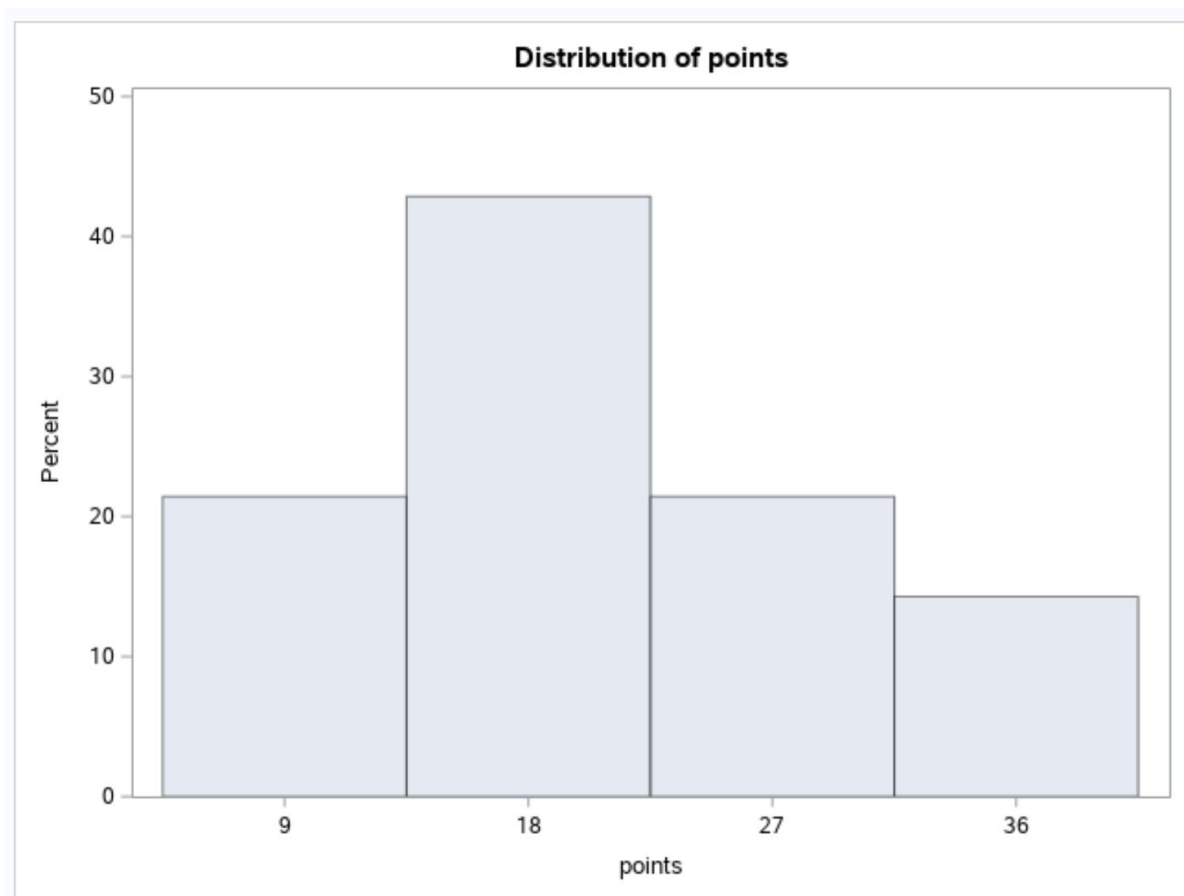
Conversely to seeking granular detail, there are frequent analytical scenarios where a broader, more consolidated view of the data distribution is far more appropriate or preferred. This preference often arises during the initial exploratory phase or when generating simplified summaries for non-technical audiences. To achieve this necessary simplification, we must strategically decrease the overall number of [bins](#) by purposefully specifying a much larger `interval` value for the `midpoints` option. This approach effectively consolidates the numerous data points into fewer, substantially wider bins, which inherently helps to smooth out minor, insignificant fluctuations and clearly highlight only the major, overarching trends.

Let us adjust our [SAS](#) code block one final time. In this iteration, we will specify `midpoints` that occur at fixed intervals of 9 units, ensuring the command still covers the established range from 9 to 36. This specific configuration will drastically reduce the bin count compared to both the automated default and the finely-tuned histogram (`Interval=3`), offering the highest-level summary of the `points` distribution possible while retaining structure.

```
/*create histogram for points variable with custom bins (Interval=9)*/  
proc univariate data=my_data;
```

```
histogram points / midpoints=(9 to 36 by 9);  
run;
```

As clearly depicted in the histogram below, the visual representation now features considerably fewer [bins](#). This intentional, broader binning strategy places great emphasis on the overall macro-shape of the distribution, making it considerably easier for the viewer to rapidly identify the primary central clusters or assess the general symmetry of the data, without being distracted by minor, minute details. Such a simplified view is frequently the preferred format for executive summaries, preliminary reporting, or initial exploratory phases where a highly simplified, yet accurate, understanding of the data is paramount for decision-making.



## Choosing the Optimal Bin Width for Your Histogram

The critical decision regarding the ideal number of [bins](#), or equivalently, the optimal bin width, to utilize in a [histogram](#) is rarely a straightforward, formulaic calculation. Instead, it often demands a judicious blend of expert subjective judgment and rigorous statistical considerations. It is important to acknowledge that there exists no singular "correct" or universally optimal number of bins; the

best choice is heavily dependent on several internal factors: the intrinsic nature and characteristics of your specific [data](#), the total sample size or number of [observations](#) available, and, most importantly, the specific insights you are aiming to extract or effectively communicate through the final visualization.

To assist analysts in this complex determination, several statistical rules of thumb have been developed over the years to suggest an appropriate number of bins. Prominent examples include the widely used Sturges' Rule, the robust Freedman-Diaconis Rule, and Scott's Rule. These mathematical methods systematically incorporate factors such as the full range of the data and the sample size to propose a statistically suitable bin count. While these established rules can indeed furnish a highly valuable starting point for the analysis, they must not be followed rigidly or blindly. They function best as powerful initial guidance, but they can never fully replace the necessity of careful visual examination and critical, expert evaluation of the resulting histogram output.

Ultimately, the process of selecting the most analytically effective bin width is inherently an iterative and exploratory endeavor. It requires the analyst to actively experiment with various ``interval`` values within the ``midpoints`` statement, systematically generating and reviewing multiple [histograms](#). The final selection should be based on a visual assessment of which representation most clearly and effectively reveals the true underlying distribution characteristics of your [variable](#). When making this choice, always reflect critically on the narrative your data is attempting to convey and how different binning strategies either successfully emphasize or unintentionally de-emphasize various crucial aspects of that data narrative.

## Conclusion: Mastering Histogram Bins in SAS

The technical capability to precisely and deliberately control the count and the width of [bins](#) when working with a [SAS histogram](#) is not merely a technical option--it is a foundational skill set for achieving truly effective [data visualization](#) and advanced statistical analysis. By expertly utilizing the ``midpoints`` statement, which is housed within the powerful ``PROC UNIVARIATE`` procedure, analysts gain the necessary leverage to move decisively beyond generalized default settings and instead engineer visualizations that accurately and eloquently reflect the subtle nuances embedded within their [data](#).

As vividly demonstrated through our practical example involving the basketball player dataset, adjusting the ``interval`` parameter of the ``midpoints`` option provides immediate, tangible results. It allows the analyst to fluidly choose between two extremes: increasing the level of detail for a highly granular inspection, or reducing the detail for a concise, high-level overview summary. This intrinsic flexibility ensures that your resulting [histograms](#) are always optimally configured for the specific analytical question being addressed, invariably leading to stronger, more insightful, and defensible conclusions.

We strongly encourage all users to dedicate time to experimenting with various **`midpoints`** configurations within their own [SAS](#) programming projects. Through systematic practice and visual feedback, you will rapidly develop the essential intuitive understanding of precisely how diverse binning strategies impact the interpretation of your data distributions, ultimately culminating in the production of more robust and impactful statistical analyses. For continued professional development, we recommend reviewing the official SAS documentation regarding the **`UNIVARIATE Procedure`** and exploring supplementary tutorials focused on advanced data visualization techniques available within the SAS system.