

Learning to Select the First N Rows of a Dataset in SAS

Authored by
Mohammed looti

October 31, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Learning to Select the First N Rows of a Dataset in SAS*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=7342>

Efficiently managing and analyzing large [datasets](#) is a core responsibility of any professional using [SAS](#) programming. Data analysts frequently need to isolate a small portion of the data, particularly the initial observations, for crucial tasks such as debugging code, performing rapid data validation checks, or focusing specific analyses on the most recent entries. This comprehensive article serves as an expert guide, detailing the two most reliable and commonly employed methods for precisely selecting the first 'N' [rows](#) from a SAS dataset. We provide clear, step-by-step explanations and practical coding examples to ensure mastery of these techniques.

Gaining proficiency in these subsetting techniques is fundamental for optimizing your data manipulation workflow within the SAS environment. We will explore the essential programming constructs that facilitate this process, ensuring you understand not only the syntax required for implementation but also the underlying logical mechanism. Both methods are built upon the robust framework of the DATA step, which is the cornerstone of SAS for creating, modifying, and transforming datasets.

The Foundation: Understanding the DATA Step and Automatic Variables

The DATA step is the foundational component in SAS used for all data manipulation, transformation, and creation of new datasets. It operates by processing data observation by observation, allowing programmers to implement intricate logical and conditional operations. When initiating a process that involves reading existing data, the [SET statement](#) is indispensable, as it explicitly instructs SAS to retrieve observations sequentially from one or more designated SAS datasets.

A pivotal element for selecting specific records based on their sequential position is the [_N_ automatic variable](#). This special variable is automatically generated and maintained by SAS throughout the execution of the DATA step, and it reliably represents the current iteration count. Essentially, `_N_` records the number of times the DATA step has begun processing a new observation. Consequently, for the first observation read, `_N_` equals 1; for the second, `_N_` equals 2, and so forth. This inherent counting mechanism makes `_N_` perfectly suited for precisely selecting rows based on their sequential order in the input data.

To exert precise control over which observations are retained and written to the newly created dataset, we utilize the [IF statement](#) in conjunction with the [OUTPUT statement](#). The IF statement allows for the conditional execution of SAS statements, enabling us to specify exact criteria for data inclusion. The OUTPUT statement explicitly writes the current observation (the one currently being processed) to the dataset under construction. While SAS typically outputs observations implicitly at the conclusion of each DATA step iteration, using an explicit OUTPUT statement within an IF condition grants the programmer granular control, ensuring only observations that meet the specified positional criteria are included in the final output dataset.

Method 1: Precision Subsetting for the Initial Observation

When the analytical goal is restricted to isolating only the very first observation from a large dataset, Method 1 offers the most direct and resource-efficient solution. This technique skillfully leverages the fundamental nature of the `_N_` automatic variable, which tracks the iteration count of the DATA step. By establishing a conditional check--specifically, determining if the value of `_N_` is precisely equal to 1--we can accurately pinpoint and select the singular first row processed by SAS.

The syntax for this method is highly streamlined, comprising a DATA step that reads data via the SET statement from your original dataset, and an IF statement that conditions the execution of the OUTPUT statement. This combination guarantees that only the initial record satisfying the condition is written to the new dataset. A critical aspect of this technique is that once the first observation has been outputted, the DATA step continues to execute for the remainder of the input dataset, but no subsequent outputs will occur because the condition `_N_ = 1` will no longer be true for any following observation.

```
data first_row;  
set original_data;  
if _N_ = 1 then output;  
run;
```

Examining the code above reveals the implementation logic:

`data first_row;` initiates the creation of a new dataset named `first_row` using a DATA step.

`set original_data;` instructs SAS to read observations sequentially from the existing input dataset `original_data`.

`if _N_ = 1 then output;` represents the core conditional logic. It checks if the current observation being processed corresponds to the first iteration (where `_N_` equals 1). If this is true, the OUTPUT statement explicitly writes this observation to the `first_row` dataset. Because an explicit OUTPUT is used, SAS will only write records when this specific condition is met, effectively stopping after the first record.

`run;` executes the defined DATA step, creating the final subset.

Preparing the Environment: Creating the Sample Dataset

To provide a tangible demonstration of these subsetting methods, we will first construct a simple, internal [dataset](#). This example dataset simulates hypothetical basketball team statistics, including

variables for team name, points scored, and total rebounds. Utilizing an in-stream sample allows us to clearly illustrate the functionality of the SAS code without relying on external file dependencies, ensuring reproducibility and clarity.

The following SAS code block executes two primary functions. First, it creates the `original_data` dataset. We employ the [DATALINES statement](#) to input the data directly within the program, followed by the INPUT statement which defines the structure: `team` is defined as a character [variable](#) (denoted by the \$ sign), and `points` and `rebounds` are defined as standard numeric variables. Second, immediately after the dataset creation, we use [PROC PRINT](#) to display the contents of `original_data`, thereby confirming that the sample data has been loaded and structured correctly before we proceed with the subsetting demonstrations.

```
/*create dataset*/  
data original_data;  
input team $ points rebounds;  
datalines;  
Warriors 25 8  
Wizards 18 12  
Rockets 22 6  
Celtics 24 11  
Thunder 27 14  
Spurs 33 19  
Nets 31 20  
Mavericks 34 10  
Kings 22 11  
Pelicans 39 23  
;  
run;  
  
/*view dataset*/  
proc print data=original_data;  
run;
```

Obs	team	points	rebounds
1	Warriors	25	8
2	Wizards	18	12
3	Rockets	22	6
4	Celtics	24	11
5	Thunder	27	14
6	Spurs	33	19
7	Nets	31	20
8	Maverick	34	10
9	Kings	22	11
10	Pelicans	39	23

Implementation Showcase: Extracting the First Row

With our `original_data` dataset successfully created and verified, we can now execute Method 1 to extract only the single initial observation. This technique proves exceptionally valuable when the analyst needs to quickly inspect the dataset's starting structure, confirm initial data integrity, or select a single, representative record for prototype testing purposes. The code presented below executes this process, resulting in a new dataset that contains solely the statistics for the first team listed in our sample data.

The DATA step named `first_row` is executed, systematically reading `original_data` observation by observation. The IF condition `_N_ = 1` evaluates to true only during the processing of the first record. When this condition is met, the OUTPUT statement writes that specific observation to the new dataset `first_row`. Critically, for all subsequent observations, the condition fails, ensuring that no further rows are written, thereby achieving the desired single-row extraction. Subsequently, [PROC PRINT](#) is utilized to display the contents of the newly created `first_row` dataset, allowing us to immediately verify the successful outcome.

```
/*create new dataset that contains only the first row*/
```

```
data first_row;  
set original_data;  
if _N_ = 1 then output;  
run;
```

```
/*view new dataset*/  
proc print data=first_row;  
run;
```

Obs	team	points	rebounds
1	Warriors	25	8

As the resulting output clearly illustrates, the `first_row` dataset accurately contains only the initial observation ("Warriors") from our original source data. This confirms the efficacy and precision of using the `_N_ = 1` condition for single-row extraction purposes.

Method 2: Flexible Selection of the First N Records

When the requirement is to select a defined number of initial observations--for example, the first 5, 10, or any arbitrary 'N' rows--Method 2 provides the necessary flexibility and control. This technique operates on the same principle as Method 1, relying heavily on the `_N_` automatic variable; however, it incorporates a crucial modification to the conditional logic. Instead of checking for strict equality to 1, we implement a check to see if the value of `_N_` is less than or equal to the desired number 'N'.

This approach is highly adaptable, allowing programmers to control the exact size of the subset simply by adjusting the value 'N' within the [IF statement](#). This makes Method 2 an invaluable tool for generating smaller, highly manageable subsets of extremely large datasets, which are ideal for preliminary data exploration, quick debugging cycles, or focused analysis on time-series data where the most recent entries are prioritized. Once the condition `_N_ <= N` is no longer true (i.e., `_N_` is greater than N), the DATA step will continue to process the remaining observations from the input dataset but will cease writing any more records to the new dataset, thereby preserving computational efficiency.

```
data first_N_rows;
set original_data;
if _N_ <= 5 then output; /*select first 5 rows*/
run;
```

A detailed breakdown of the code for selecting the first five rows:

`data first_N_rows;` initializes the DATA step, designating the name of the new output dataset as `first_N_rows`.

`set original_data;` specifies that the processing should sequentially read data from `original_data`.

`if _N_ <= 5 then output;` contains the core conditional logic. It checks if the current observation number (`_N_`) is less than or equal to 5. If this condition is satisfied, the `OUTPUT` statement explicitly writes the observation to `first_N_rows`.

`run;` executes the `DATA` step, creating the final subset containing the first five rows.

Implementation Showcase: Selecting a Specific Number of Rows

We now proceed to apply Method 2 to our `original_data` dataset, demonstrating the extraction of the first five observations. This practical implementation highlights how easily large datasets can be condensed into smaller, targeted subsets, a process invaluable for focused analysis or rapid prototyping.

The following SAS code executes a `DATA` step to create `first_N_rows`, reading `original_data` in sequence. The conditional statement `if _N_ <= 5 then output;` acts as the filter, ensuring that only records where the automatic variable `_N_` is 5 or less are written to the new dataset, thus precisely capturing the first five rows. Once `_N_` surpasses 5, the condition evaluates to false, and subsequent observations are automatically excluded from the output. Following the subsetting `DATA` step, another [PROC PRINT](#) statement is utilized to display the contents of `first_N_rows`, visually confirming the successful creation of the five-row subset.

```
/*create new dataset that contains first 5 rows of original dataset*/
```

```
data first_N_rows;
```

```
set original_data;
```

```
if _N_ <= 5 then output;
```

```
run;
```

```
/*view new dataset*/
```

```
proc print data=first_N_rows;
```

```
run;
```

Obs	team	points	rebounds
1	Warriors	25	8
2	Wizards	18	12
3	Rockets	22	6
4	Celtics	24	11
5	Thunder	27	14

The resulting dataset, `first_N_rows`, distinctly presents the first five observations from our source data, exactly as intended by the programming logic. This successful outcome underscores the flexibility achieved by combining the `_N_` automatic variable with a less-than-or-equal-to condition, providing a powerful mechanism to extract any desired number of leading records. To subset a different quantity of starting rows, analysts simply need to modify the value following the `_N_ <=` operator in the code block above to match their required 'N' value.

Summary of Techniques and Best Practices

Acquiring proficiency in selecting the first 'N' rows of a dataset is a foundational and indispensable skill in professional SAS programming, granting analysts significant efficiency and robust control over crucial data manipulation tasks. Regardless of whether the requirement is to extract a single observation for immediate inspection or a larger subset for focused, iterative analysis, the methods detailed herein--which rely on leveraging the `_N_` automatic variable within a DATA step structure--provide straightforward, reliable, and highly scalable solutions.

These techniques are not only powerful in their simplicity but are also highly flexible across various analytical contexts. By developing a deep understanding of the core logic--how SAS sequentially processes data observation by observation, and critically, how to conditionally control the OUTPUT statement--you can readily extend these principles to address a broad spectrum of other complex data subsetting and transformation challenges. We strongly recommend that users practice these methods extensively using their own production or test datasets to thoroughly cement their understanding and explore the full potential of these efficient programming constructs.

Additional Resources

The following tutorials provide further explanations on how to perform other essential and common data manipulation tasks in SAS:

[How to Filter Rows in SAS](#)

[How to Sort Data in SAS](#)

[How to Merge Datasets in SAS](#)