

Learning Data Subsetting Techniques in SAS: A Practical Guide with Examples

Authored by
Mohammed Iooti

October 31, 2025

RECOMMENDED CITATION

Mohammed Iooti (2025). *Learning Data Subsetting Techniques in SAS: A Practical Guide with Examples*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=7367>

Introduction: Mastering Data Subsetting in SAS

In the expansive world of data management and advanced analytics, the capability to efficiently refine and focus large datasets is absolutely paramount. One of the most fundamental operations required by data professionals is [subsetting](#), which is the selective process of extracting a specific portion of data based on predefined criteria. This technique is not merely about convenience; it is a critical step for optimizing processing speed, managing memory resources, and ensuring that subsequent statistical analyses are performed only on the most relevant information.

[SAS](#), or the **S**tatistical **A**nalysis **S**ystem, stands as a premier software suite relied upon globally for tasks ranging from business intelligence to sophisticated predictive modeling. Within the [SAS](#) environment, users are provided with a variety of intuitive and powerful commands to perform precise data subsetting. These commands allow the user to control exactly which observations (rows) and [variables](#) (columns) are transferred or excluded when creating a new [dataset](#) from an existing one.

This comprehensive guide is dedicated to exploring the three most effective and frequently employed methods for subsetting a [dataset](#) in [SAS](#). We will delve into the specific techniques used for selecting desired columns, excluding irrelevant columns, and filtering rows based on complex logical conditions. Each method will be presented with clear, formal explanations and practical, runnable code examples designed to illuminate their application in real-world scenarios.

Fundamental Approaches to Subsetting Data in SAS

The cornerstone of data preparation and manipulation in [SAS](#) programming is the [DATA step](#). This versatile procedural block is essential for reading raw data, creating new datasets, modifying existing data structures, and executing intricate transformations. Consequently, nearly all robust subsetting operations are performed directly within a [DATA step](#), where the user defines the source data and specifies the precise rules governing the construction of the desired output [dataset](#).

A solid understanding of these fundamental approaches is absolutely necessary for anyone seeking mastery over data management in [SAS](#). Each technique we will examine offers a distinct strategy for tailoring the data to meet specific analytical requirements. Whether the objective is to streamline a voluminous [dataset](#) by isolating only a few key [variables](#) or to dramatically narrow the number of observations to those that satisfy strict logical conditions, [SAS](#) provides elegant and high-performance solutions.

The following list outlines the three primary methods for [subsetting](#) data in [SAS](#), providing a concise preview of the code structure used for each operation. These examples demonstrate whether the focus is on selecting columns, excluding columns, or filtering rows.

Method 1: Column Selection using the [KEEP statement](#)

This approach is optimal when the relevant [variables](#) are few and clearly identified. The [KEEP statement](#) explicitly defines the columns that must be retained in the resulting [dataset](#). All variables that are not listed in the [KEEP statement](#) are automatically and permanently excluded from the output.

```
data new_data;  
set original_data;  
keep var1 var3;  
run;
```

Method 2: Column Exclusion using the [DROP statement](#)

Conversely, the [DROP statement](#) is best suited for scenarios where the original data contains a large number of [variables](#), but only a small subset needs to be removed. By listing the [variables](#) to be dropped, all unmentioned variables are automatically preserved. This often proves to be a more concise and efficient coding solution than using the [KEEP statement](#) when the goal is to exclude only a few columns.

```
data new_data;  
set original_data;  
drop var4;  
run;
```

Method 3: Conditional Row Filtering using the [IF...THEN DELETE statement](#)

This powerful method allows for the selection of observations (rows) based on one or more specific criteria. The [IF...THEN DELETE statement](#) is the cornerstone of conditional [subsetting](#) in [SAS](#). It evaluates a logical expression for every row; if the condition is met (evaluates to true), that observation is immediately deleted and is not written to the output [dataset](#). This effectively serves as the primary mechanism for filtering data.

```
data new_data;  
set original_data;  
if var1 < 25 then delete;  
run;
```

Setting the Stage: Our Sample Dataset

To provide practical, observable results for these [subsetting](#) techniques, we must first establish a simple, representative [dataset](#). This source data, which we will name `original_data`, will contain basic metrics for several fictional basketball teams, including their scoring performance and rebounding statistics. This dataset will serve as the consistent baseline for all subsequent examples, allowing us to clearly track the effect of each specific subsetting operation.

The structure of the `original_data` [dataset](#) comprises three [variables](#): `team` (a character variable identifying the team name), `points` (a numeric variable representing points scored), and `rebounds` (a numeric variable indicating rebounds collected). This clear, focused structure allows us to demonstrate both column-based selection/exclusion and row-based conditional filtering effectively.

The following [SAS](#) code block illustrates the creation of this sample [dataset](#). The creation process utilizes the [DATA step](#) in conjunction with the `INPUT` statement and `DATALINES` to input the data directly. After the dataset is created, we employ the [PROC PRINT](#) procedure to display its contents, establishing a visual reference for comparison with the forthcoming subsetted results.

```
/*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  
;  
run;  
  
/*view dataset*/  
proc print data=original_data;
```

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

Practical Application 1: Subsetting by Column Selection (The KEEP Statement)

When the primary analytical requirement is to focus exclusively on a limited selection of the available [variables](#), the [KEEP statement](#) provides the most direct and resource-efficient methodology. By explicitly naming the variables intended for retention, [SAS](#) constructs a new [dataset](#) that contains only those specified columns, ensuring all other columns are effectively discarded. This technique is invaluable for dataset simplification, minimizing the computational load, and preparing data for highly focused statistical procedures.

Let us apply this principle to our `original_data` [dataset](#), which currently holds `team`, `points`, and `rebounds`. If our current analysis is solely concerned with team identification and their scores, we must use the [KEEP statement](#) to select only the `team` and `points` [variables](#). The structure within the [DATA step](#) remains consistent: we declare a new [dataset](#) (`new_data`), specify the input (`original_data`), and then implement our [subsetting](#) command.

The [SAS](#) code provided below demonstrates how to subset the `original_data` [dataset](#), preserving only the `team` and `points` columns. The [KEEP statement](#) ensures that only these two variables are written to the `new_data` dataset. We conclude the process by using [PROC PRINT](#) to visualize the resulting structure.

```
/*create new dataset*/  
data new_data;  
set original_data;  
keep team points;  
run;  
  
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points
1	Warriors	25
2	Wizards	18
3	Rockets	22
4	Celtics	24
5	Thunder	27
6	Spurs	33
7	Nets	31

As clearly illustrated by the output, the `new_data` [dataset](#) now exclusively features the `team` and `points` variables, while the `rebounds` variable has been successfully eliminated. This outcome decisively demonstrates the precision and effectiveness of the [KEEP statement](#) in defining the exact column structure of the output dataset.

Practical Application 2: Subsetting by Column Exclusion (The DROP Statement)

In contrast to the selective inclusion offered by the [KEEP statement](#), the [DROP statement](#) provides an equally effective but inverted method for managing [variables](#). This methodology is particularly advantageous when dealing with datasets that possess a vast number of variables, yet the user only needs to remove a small, identifiable subset. Instead of listing dozens of variables to keep, the user simply specifies those to discard, and all remaining variables are automatically carried forward into the new dataset.

Revisiting our `original_data` [dataset](#), let us imagine that our analysis requires all information except for the `points` variable. In this specific scenario, the [DROP statement](#) presents the most concise coding option. We simply insert the command `DROP points;` within our [DATA step](#), and [SAS](#) subsequently generates a new [dataset](#) containing only `team` and `rebounds`.

The code below demonstrates the application of the [DROP statement](#). Similar to the [KEEP statement](#), this operation is encapsulated within a [DATA step](#), resulting in the new `new_data` [dataset](#). We utilize [PROC PRINT](#) immediately afterward to confirm that the `points` column has been successfully excluded from the output.

```
/*create new dataset*/  
data new_data;  
set original_data;  
drop points;
```

```
run;
```

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

```
proc print data=new_data;
```

Obs	team	rebounds
1	Warriors	8
2	Wizards	12
3	Rockets	6
4	Celtics	11
5	Thunder	14
6	Spurs	19
7	Nets	20

The resultant output clearly verifies that the `new_data` [dataset](#) now contains only the `team` and `rebounds` variables. The `points` variable has been effectively excluded, demonstrating the power of the [DROP statement](#) for selective column removal. Both the [KEEP statement](#) and the [DROP statement](#) are fundamentally important for column-based [subsetting](#); the ultimate choice between them depends solely on which list of variables (those to retain or those to remove) is shorter and more convenient to specify.

Practical Application 3: Subsetting by Conditional Row Selection (IF...THEN DELETE)

Moving beyond the selection or exclusion of entire columns, data [subsetting](#) frequently requires filtering observations (rows) based on specific, dynamic criteria. For this purpose, the [IF...THEN DELETE statement](#), executed within the [DATA step](#), is an exceptionally versatile and critical tool. It allows the programmer to define a condition; if any row satisfies that condition (evaluates to true), the row is immediately prevented from being written to the output [dataset](#). This mechanism ensures that only the rows that successfully fail the deletion criteria are retained.

Let us apply this principle to our `original_data` dataset. Suppose our analysis requires only those teams that achieved a strong scoring performance of 25 points or more. To achieve this selective filtering, we must instruct [SAS](#) to delete any row where the value in the `points` column is strictly less than 25. This conditional deletion ensures that only teams demonstrating the required scoring threshold are preserved for subsequent analytical steps.

The [SAS](#) code below demonstrates this conditional row selection process. The [DATA step](#)

sequentially processes each row of `original_data`. For every row where the condition `points < 25` is met, the `THEN DELETE` action is triggered, ensuring that row is excluded from the final `new_data`. Finally, [PROC PRINT](#) is used to display the filtered results.

```
/*create new dataset*/  
data new_data;  
set original_data;  
if points < 25 then delete;  
run;  
  
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points	rebounds
1	Warriors	25	8
2	Thunder	27	14
3	Spurs	33	19
4	Nets	31	20

The resulting output confirms that only the rows where `points` are 25 or greater have been successfully retained. This clearly showcases the fundamental application of the [IF...THEN DELETE statement](#) for precise row-level [subsetting](#).

Expanding Conditions: Using Logical Operators

The utility of conditional subsetting can be significantly amplified by integrating multiple conditions using [logical operators](#). These operators facilitate the creation of far more intricate filtering rules, enabling the selection of rows that must satisfy several criteria simultaneously or any one of a group of criteria. The two primary logical operators utilized in [SAS](#) are [OR](#) (represented by the pipe symbol `|` or the keyword `OR`) and [AND](#) (represented by the ampersand `&` or the keyword `AND`).

The [OR operator](#) (`|`) dictates that a row should be deleted if **any** of the specified conditions evaluate to true. For example, if we seek to remove rows where `points` is less than 25 **OR** `rebounds` is less than 10, a row will be deleted if either criterion (or both simultaneously) is met. This results in a more aggressively filtered [dataset](#) with fewer observations.

```
/*create new dataset*/  
data new_data;
```

```
set original_data;  
if points < 25 | rebounds < 10 then delete;  
run;
```

```
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points	rebounds
1	Thunder	27	14
2	Spurs	33	19
3	Nets	31	20

The output above clearly demonstrates the effect of the [OR operator](#). Rows are deleted if either the points condition (e.g., Wizards, Celtics) or the rebounds condition (e.g., Warriors, Rockets) is satisfied, resulting in a dataset with only three retained observations.

Conversely, the [AND operator](#) (&) is deployed when a row should be deleted only if **all** specified conditions are simultaneously true. For instance, if the instruction is to remove rows where `points` is less than 25 **AND** `rebounds` is less than 10, a row will only be deleted if both conditions are met for that specific observation. This approach generates a more selectively filtered [dataset](#), typically retaining more observations than a comparable [OR](#) condition.

```
/*create new dataset*/  
data new_data;  
set original_data;  
if points < 25 & rebounds < 10 then delete;  
run;
```

```
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points	rebounds
1	Warriors	25	8
2	Wizards	18	12
3	Celtics	24	11
4	Thunder	27	14
5	Spurs	33	19
6	Nets	31	20

Observing the output derived from the [AND operator](#) reveals that only the 'Rockets' row was deleted. This is because it was the only team where both `points` were less than 25 (22) AND `rebounds` were less than 10 (6). Other rows, such as 'Warriors' (points=25, rebounds=8), satisfy only one condition and are therefore retained. This highlights the precise, highly controlled filtering capability offered by logical operators in complex subsetting tasks.

Conclusion and Further Exploration

The ability to efficiently execute data [subsetting](#) operations is an indispensable cornerstone skill for every [SAS](#) programmer and analyst. This guide has systematically presented the three most powerful and common methods for refining datasets: defining specific columns to retain using the [KEEP statement](#), excluding unnecessary columns with the [DROP statement](#), and filtering observations based on complex, conditional logic using the [IF...THEN DELETE statement](#) combined with [OR](#) and [AND](#) operators.

These core techniques, primarily executed within the [DATA step](#), offer the necessary flexibility and precision to prepare your data for virtually any analytical challenge. By effectively managing and [subsetting](#) your data, you can significantly enhance the performance of your [SAS](#) programs, sharpen the focus of your analyses, and guarantee the integrity of your results by working exclusively with the most relevant portions of your source [dataset](#).

We strongly encourage you to apply and experiment with these methods across your own data management workflows. Developing proficiency in these efficient subsetting practices is a foundational skill that will profoundly contribute to your overall expertise in [SAS](#) programming and data handling.

Additional Resources

To continue building upon your [SAS](#) expertise, consider exploring related advanced data manipulation techniques and other common data preparation tasks. The following resources

provide valuable guidance on topics that complement the [subsetting](#) skills demonstrated in this tutorial:

SAS Official Documentation: For comprehensive and authoritative details on any [SAS](#) statement, procedure, or function.

Online [SAS](#) Communities: An excellent platform for seeking assistance, asking specific coding questions, and finding solutions shared by experienced users.

Books and Courses on [SAS](#) Programming: Recommended for pursuing structured learning paths and gaining in-depth knowledge of statistical computing.