

Learning SAS: Filtering Data with the WHERE Option and SET Statement

Authored by
Mohammed loot

November 14, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning SAS: Filtering Data with the WHERE Option and SET Statement*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=1489>

In the demanding environment of statistical analysis and business intelligence, [SAS](#) stands out as an essential tool for managing, transforming, and analyzing immense volumes of information. Achieving efficiency in data processing often hinges on mastering core techniques, one of the most vital being the integration of the [WHERE option](#) with the [SET statement](#). This powerful combination empowers developers and analysts to execute dynamic, high-speed filtering of input data right at the source. The primary objective of this technique is to create a new, highly focused [dataset](#) by extracting only those observations from the source that strictly adhere to specific, predefined [conditions](#). This methodology is indispensable for effective [data manipulation](#), ensuring that all subsequent analytical procedures are performed exclusively on the most relevant subset of data, which dramatically improves performance and clarity in results.

The capacity to embed filtering logic directly into the data input mechanism--as opposed to applying filters after the entire source dataset has been read into memory--provides a massive performance advantage. This benefit is particularly pronounced when working with enterprise-scale data warehouses where minimizing I/O operations is critical. Throughout this comprehensive article, we will meticulously explore the practical applications and technical implementation details of using the **WHERE option** alongside the **SET statement** in [SAS](#) programming. We will break down two fundamental methodologies: filtering based on a singular, straightforward criterion, and executing complex filtering tasks using intricate logical expressions that involve multiple, interwoven conditions. By mastering these approaches, users can significantly optimize their data preparation workflow, guaranteeing that analyses are conducted swiftly and precisely on the required data subset.

Understanding the Efficiency of WHERE in the SET Statement

The [SET statement](#) is a foundational element within the [DATA step](#) in [SAS](#), designed to sequentially read observations from an existing source [SAS dataset](#). When the [WHERE option](#) is strategically appended to the source dataset name within the parentheses of the `SET` statement, it activates an immediate and highly efficient mechanism for data filtering. The core advantage lies in the timing: this filtering occurs as the observations are being streamed and read from the source dataset, but critically, before they are processed or written to the target dataset being created. This "on-the-fly" processing is inherently superior in efficiency compared to alternatives like reading all records into memory first and then applying a separate `WHERE` statement outside the `SET` option, or relying on complex `IF-THEN` conditional logic, especially when dealing with massive datasets where resource utilization is a major concern.

To implement this, the standard syntax dictates placing the **WHERE option** keyword directly inside the parentheses following the source dataset name. This keyword is then followed by a logical expression, which precisely defines the required inclusion criteria. These criteria are formalized as [conditions](#) that must evaluate to true for an observation to be selected and included in the resultant

dataset. The inherent flexibility and conciseness of this structure allow for highly readable and maintainable [SAS](#) code, often eliminating the need for more complex procedural steps, such as utilizing `PROC SQL`, solely for the purpose of subsetting data. Mastering this specific syntax is paramount for writing streamlined, performant data management routines, thereby establishing a solid foundation for advanced [data manipulation](#) tasks.

Method 1: Subsetting Data Using a Single Criterion

The most straightforward, yet immensely practical, use case for the [WHERE option](#) is its application within the [SET statement](#) to filter a source [dataset](#) based on just one, clearly defined criterion. This technique is invaluable when the analytical goal is to isolate rows where a specific variable satisfies a simple logical test. Examples of such tests include evaluating whether a numerical variable exceeds a defined threshold, determining if a character variable exactly matches a specified string value, or checking if a value falls within a predefined range using operators such as `BETWEEN`.

The following [SAS](#) code snippet clearly illustrates the fundamental structure necessary for applying a single [condition](#) during the data reading phase. This structure facilitates the creation of a new dataset that contains only those observations from the source that perfectly align with the single specified criterion, thereby achieving an immediate reduction in the size and complexity of the data requiring further analysis.

```
data new_data;  
set my_data (where = (points>20));  
run;
```

In this specific, illustrative example, a new dataset named **new_data** is instantiated. This dataset strictly inherits observations from the source dataset, **my_data**, but only if the recorded value in the **points** column is numerically greater than 20. This stringent, single-point filter guarantees that only records satisfying this precise numerical criterion are carried forward, drastically streamlining subsequent analytical stages and ensuring a focused scope for the entire investigation.

Method 2: Integrating Logical Operators for Complex Filtering

When data filtering requirements extend beyond a simple test on a single variable, the [WHERE option](#) provides the capability to integrate multiple [conditions](#) through the use of essential logical operators, most commonly [OR](#) and [AND](#). This advanced capability grants analysts granular control over their data selection process, enabling them to select observations based on intricate rules that often involve relationships between several variables simultaneously. This powerful, advanced subsetting technique is crucial for generating highly specific data cohorts for detailed study or

customized reporting, making complex [data manipulation](#) tasks both efficient and straightforward within the [DATA step](#).

The following code snippet demonstrates the practical application of multiple conditions utilizing the inclusive [OR operator](#). This logical structure dictates an inclusive rule: an observation will be included in the new dataset if it satisfies at least one of the specified conditions. This inclusive logic is perfectly suited for combining disparate groups that share a common analytical interest, effectively creating a union of records that meet varied criteria.

```
data new_data;  
set my_data (where = (points>20 or team="Rockets"));  
run;
```

In this refined scenario, the resulting dataset, **new_data**, is populated by drawing observations from **my_data** where either the **points** variable exceeds 20 or the **team** variable is exactly equal to the string "Rockets". This illustration highlights the versatility of combining different data types and comparison criteria, enabling the achievement of highly nuanced filtering outcomes. Conversely, employing the [AND operator](#) mandates a stricter, intersectional requirement: an observation must simultaneously satisfy **all** specified conditions to be selected. Understanding the critical difference between **OR** (union) and **AND** (intersection) logic is absolutely fundamental to accurately defining and generating the precise data subsets required for any comprehensive analysis in [SAS](#).

Practical Implementation: Defining the Source Data

To provide a concrete and verifiable foundation for demonstrating these two powerful filtering methods, it is essential that we first define and populate a sample [dataset](#) within the [SAS](#) environment. This foundational dataset, which we will consistently refer to as **my_data**, contains key performance information related to various sports teams, specifically tracking metrics such as total points scored and total assists recorded. This consistent source data will serve as the reliable basis for all subsequent filtering and data subsetting examples, thereby ensuring maximum clarity, repeatability, and ease of understanding in our demonstrations.

The following [DATA step](#) code leverages the ``INPUT`` and ``DATALINES`` statements to quickly construct and populate the **my_data** dataset. Immediately upon the successful creation of the dataset, the ``PROC PRINT`` statement is invoked. This procedure displays the raw, unfiltered contents of this newly formed dataset. This initial visual inspection is a crucial step, as it provides a clear baseline of the complete data against which we can accurately measure the success, precision, and efficiency of our filtering operations in the subsequent sections of this demonstration.

```
/*create dataset*/
data my_data;
input team $ points assists;
datalines;
Mavs 22 10
Rockets 12 14
Spurs 29 8
Kings 13 10
Warriors 44 10
Heat 18 8
Magic 11 5
Pelicans 19 3
Blazers 12 8
;
run;

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

Obs	team	points	assists
1	Mavs	22	10
2	Rockets	12	14
3	Spurs	29	8
4	Kings	13	10
5	Warriors	44	10
6	Heat	18	8
7	Magic	11	5
8	Pelicans	19	3
9	Blazers	12	8

As clearly confirmed by the printed output displayed above, the **my_data** dataset has been successfully initialized and populated with nine distinct observations. Each observation represents a complete record for a team, coupled with its associated points and assists data. This established source dataset now provides the ideal context for showcasing the enhanced efficiency and precision afforded by integrating the [WHERE option](#) directly into the data reading phase of the [SET statement](#).

Demonstration of Single-Condition Filtering in Practice

With the source data, **my_data**, successfully created, we can now proceed to apply the first primary filtering technique: subsetting the data based purely on a single, numerical [condition](#). Our specific objective in this demonstration is to construct a new dataset, which we designate **new_data**, that exclusively contains records for those teams whose cumulative **points** score strictly exceeds the numerical threshold of 20. This action serves to clearly illustrate how to efficiently target and isolate specific numerical cohorts using streamlined and highly performant SAS syntax.

The subsequent [SAS](#) code precisely implements this filter by embedding the **WHERE option** within the parentheses of the source dataset reference in the **SET statement**. Immediately following the successful generation of the filtered dataset **new_data**, we utilize `PROC PRINT` to visually confirm the results. This provides a direct, side-by-side comparison against the initial, unfiltered data displayed earlier, allowing us to immediately verify the filtering accuracy.

```
/*create new dataset*/  
data new_data;  
set my_data (where = (points>20));  
run;  
  
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points	assists
1	Mavs	22	10
2	Spurs	29	8
3	Warriors	44	10

A thorough review of the resulting output confirms the precise and successful application of the single condition filter. The **new_data** dataset now consists of exactly three observations: Mavs, Spurs, and Warriors. These are the only records selected from the original **my_data** that possessed a value in the **points** column strictly greater than 20. This conclusive outcome validates the exceptional effectiveness of the **WHERE option** in performing precise, resource-efficient data subsetting tasks based on a singular, well-defined criterion.

Demonstration of Multiple-Condition Filtering with Logical Operators

We now advance to demonstrating more sophisticated and flexible filtering logic by incorporating the use of logical operators to combine multiple conditions within the **WHERE option**. For this specific example, our objective is to populate **new_data** with observations that satisfy either of two distinct criteria: (1) the **points** score is greater than 20, **OR** (2) the **team** name is precisely the character string "Rockets". This scenario is intentionally designed to showcase the power and utility of the inclusive logical operator in handling complex, heterogeneous data selection tasks.

The syntax provided below seamlessly integrates the **WHERE option** within the **SET statement**, specifically employing the **OR operator** to link the numerical comparison condition with the character comparison condition. Following the execution of this **DATA step**, `PROC PRINT` is immediately utilized to display the filtered results, allowing us to verify that the complex logical operation was performed accurately and according to the desired rules.

```
/*create new dataset*/  
data new_data;  
set my_data (where = (points>20 or team="Rockets"));  
run;  
  
/*view new dataset*/  
proc print data=new_data;
```

Obs	team	points	assists
1	Mavs	22	10
2	Rockets	12	14
3	Spurs	29	8
4	Warriors	44	10

Upon examining the final output, we confidently observe that **new_data** now contains four observations: Mavs, Rockets, Spurs, and Warriors. This outcome perfectly confirms the behavior of the **OR operator** - an observation is selected if its **points** are greater than 20 (Mavs, Spurs, Warriors) **or** if its **team** is "Rockets" (Rockets, which is included despite having only 12 points). This flexibility is absolutely invaluable for dynamic and condition-driven data preparation processes.

It is paramount for analysts to recognize the contrasting, restrictive effect of using the **AND operator** in place of **OR**. If the expression in the code were altered to `(points>20 and`

team="Rockets")`, the resulting dataset would be entirely empty, as no single observation in our source data satisfies both stringent conditions simultaneously. Therefore, the strategic choice of logical operator--whether **AND** for required intersection or **OR** for inclusive union--is the key determinant in accurately defining the final data subset required for analysis.

Conclusion and Resources for Advanced Learning

The mastery of employing the **WHERE option** directly within the **SET statement** represents an essential and fundamental skill set for any proficient **SAS** user. This technique dramatically enhances efficiency by performing preliminary filtering precisely at the point of data input, effectively minimizing unnecessary processing and overall resource utilization. Whether the task involves applying a single, simple criterion or integrating complex, multi-variable **conditions**, this methodology guarantees precise control over the data flowing into subsequent analytical steps, ultimately leading to highly accurate and relevant insights.

By consistently applying the principles of efficient subsetting discussed throughout this article, data analysts can significantly improve both the performance and the readability of their codebases. For those committed to deepening their expertise in this powerful **programming language**, we strongly recommend exploring additional tutorials focused on various common **data manipulation** tasks and advanced procedural techniques within the expansive SAS environment.

SAS Tutorial: How to Use PROC SQL to Create a New Dataset

SAS Tutorial: How to Use the BY Statement

SAS Tutorial: How to Use the RENAME Option