

Learning How to Create New Variables in SAS: A Step-by-Step Guide

Authored by
Mohammed looti

October 31, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Learning How to Create New Variables in SAS: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=7374>

In the realm of statistical analysis and data management, the ability to generate new [variables](#) is fundamental. Whether you are standardizing scores, calculating ratios, or simply entering raw data, the [SAS](#) System provides robust and straightforward methodologies for variable creation. Understanding these mechanisms is crucial for any efficient [SAS](#) programmer seeking to transform raw information into meaningful, analysis-ready [datasets](#). We will explore the two primary approaches utilized within the [SAS](#) DATA step for introducing new variables into your workspace.

The two most common and essential techniques for generating new [variables](#) in [SAS](#) are defined by their source: either creating the variable entirely from raw input data (from scratch) or deriving the variable through computations involving existing variables within an already defined [dataset](#). Mastering both methods ensures complete flexibility in data handling and manipulation, forming the backbone of effective data preparation.

Method 1: Create Variables from Scratch

```
data original_data;  
input var1 $ var2 var3;  
datalines;  
A 12 6  
B 19 5  
C 23 4  
D 40 4  
;  
run;
```

Method 2: Create Variables from Existing Variables

```
data new_data;  
set original_data;  
new_var4 = var2 / 5;  
new_var5 = (var2 + var3) * 2;  
run;
```

The subsequent sections will delve into detailed examples illustrating how to apply each method practically. By examining the syntax and execution flow of these fundamental [SAS](#) statements, you will gain a comprehensive understanding of variable creation, enabling you to structure and manipulate your data with precision and confidence.

The Core Importance of Variable Creation in SAS Programming

Variable creation is not merely an optional step; it is the essential process of defining the analytic landscape within [SAS](#). Every piece of information, whether it is a raw measurement imported from an external file or a complex calculation derived during processing, must be stored within a defined [variable](#). The [DATA step](#) is the engine that drives this process, allowing programmers to read, transform, and output data structures. Without the ability to explicitly define new variables, sophisticated data manipulation and statistical modeling would be impossible.

When working with large or complex [datasets](#), it is common to require new variables that summarize existing data points. For instance, if you have individual scores, you might need a new [variable](#) representing the total score or the average across multiple attempts. These derived variables often simplify subsequent analysis procedures (like [PROC MEANS](#) or [PROC GLM](#)). Therefore, understanding both input methods--manual entry and derivation--is mandatory for effective programming. This foundational knowledge ensures data integrity and prepares the data for rigorous statistical interrogation.

The distinction between the two primary methods lies entirely in the origin of the data values. When creating variables from scratch, the programmer is supplying the value explicitly, usually through the code itself or via an external file definition. Conversely, when deriving variables from existing ones, the [SAS](#) compiler performs an operation (arithmetic, logical, or string manipulation) on previously defined data points. This distinction dictates which control statements--specifically [INPUT](#) or [SET](#)--must be used within the DATA step to successfully execute the variable creation logic. The following examples demonstrate how these commands are implemented in practice.

Method 1: Generating Variables from Scratch Using the DATALINES and INPUT Statements

Creating a [dataset](#) entirely from scratch is a common requirement, particularly when testing code, providing simple examples, or dealing with small amounts of manually gathered data. This method relies heavily on the [INPUT statement](#) to define the names and order of the new [variables](#), and the [DATALINES statement](#) to feed the corresponding values directly into the program. The [INPUT statement](#) establishes the structure of the data records, linking the data fields in the [DATALINES](#) section to the specified [variable](#) names.

The following detailed code demonstrates how to initialize a [dataset](#) called `original_data` containing statistical performance metrics for several sports teams. Notice that the [INPUT statement](#) dictates the three variables being created: `team`, `points`, and `rebound`. The syntax is straightforward: simply listing the desired [variable](#) names after the [INPUT](#) keyword establishes them within the [dataset](#). The [DATALINES statement](#) signals the start of the raw data entry, where

each line corresponds to a single observation, and values are matched sequentially to the variables defined in the [INPUT statement](#).

Example 1: Create Variables from Scratch

The following code shows how to create a [dataset](#) with three variables: `team`, `points`, and `rebound`:

```
/*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

It is important to notice that the values for the [variables](#) are defined directly after the [DATALINES statement](#). This approach is highly useful for creating small, self-contained examples or for inputting data that is not yet stored in an external file format. The success of this method hinges on

the sequential mapping of the data values to the variable names defined in the [INPUT statement](#), making careful attention to order crucial for correct data entry.

Understanding Data Types: Numeric vs. Character Variables

A critical consideration when defining new [variables](#) in [SAS](#) is specifying the data type. By default, [SAS](#) assumes that any new variable introduced, particularly during the [INPUT statement](#), is a numeric variable. Numeric variables are essential for mathematical operations and statistical modeling, storing integer or decimal values. If the data being read in consists solely of numbers (like `points` and `rebound` in our example), no special designation is necessary, and [SAS](#) handles them automatically as numeric.

However, when dealing with textual information--such as names, categories, or identification codes--it is mandatory to define the [character variable](#) type. To instruct [SAS](#) that a variable should store text, you must append a dollar sign (\$) immediately after the variable name in the [INPUT statement](#). In Example 1, we used `team $` to correctly define `team` as a [character variable](#), allowing it to store text strings like "Warriors" or "Celtics."

Failing to correctly identify a [character variable](#) during the input process can lead to serious data errors, as [SAS](#) will attempt to interpret non-numeric values as numeric, often resulting in missing values or incorrect data readings. This simple rule--using the dollar sign for text variables--is a cornerstone of reliable data definition and manipulation within the [DATA step](#) environment.

Method 2: Deriving New Variables from Existing Datasets Using the SET Statement

The second and arguably more frequent method of creating new [variables](#) involves computing their values based on data that already exists within a [SAS dataset](#). This is typically done for data transformation, calculating derived metrics, or applying standardization techniques. This approach utilizes the [SET statement](#), which tells [SAS](#) to read observations sequentially from a specified input [dataset](#) and copy them into the new output [dataset](#) being created.

Once the data from the original [dataset](#) is loaded into the new [DATA step](#) buffer (via the [SET statement](#)), you can define new variables simply by writing an algebraic expression. [SAS](#) evaluates this expression for each observation, automatically assigning the result to the new variable name on the left side of the equation. This makes complex data transformation remarkably intuitive, as the syntax mirrors standard arithmetic notation. For example, to calculate a variable representing half the points scored, you simply write `half_points = points / 2;`

Example 2: Create Variables from Existing Variables

The following code shows how to use the [SET statement](#) to create a new [dataset](#) (`new_data`) whose variables are created from existing [variables](#) in `original_data`:

```
/*create new dataset*/  
data new_data;  
set original_data;  
half_points = points / 2;  
avg_pts_rebs = (points + rebounds) / 2;  
run;
```

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

Obs	team	points	rebounds	half_points	avg_pts_rebs
1	Warriors	25	8	12.5	16.5
2	Wizards	18	12	9.0	15.0
3	Rockets	22	6	11.0	14.0
4	Celtics	24	11	12.0	17.5
5	Thunder	27	14	13.5	20.5
6	Spurs	33	19	16.5	26.0
7	Nets	31	20	15.5	25.5

In this example, two new [variables](#) are defined: `half_points` and `avg_pts_rebs`. The variable `half_points` simply divides the existing `points` variable by two. The variable `avg_pts_rebs` demonstrates a slightly more complex calculation, summing the `points` and `rebounds` [variables](#) before dividing by two to find the average. Note that when deriving a new variable using this method, [SAS](#) automatically determines the data type (numeric in this case) and length based on the calculation result. This flexibility makes the [SET statement](#) combined with direct assignment statements the preferred method for virtually all data preparation and transformation tasks.

Practical Applications and Advanced Transformations

While the examples provided utilize simple arithmetic, the power of variable derivation extends far beyond basic addition and division. Programmers frequently use this approach to implement conditional logic, create categorical flags, or apply complex statistical formulas. For instance, you

might use an `IF-THEN/ELSE` structure to create a binary variable (a flag) indicating whether a team's points exceed a certain threshold (e.g., `elite_scorer = (points > 30);`). Similarly, advanced mathematical functions (like `LOG`, `EXP`, or `SQRT`) can be applied directly to existing variables to normalize distributions or prepare data for specific modeling techniques.

Another crucial application involves handling missing data. A new variable can be created to identify records where critical information is absent. By using functions like `MISSING()` or checking for null values, the programmer can generate a diagnostic variable, helping analysts to manage data quality issues. Furthermore, when combining information from different sources (a process often involving the [SET statement](#) or [INPUT statement](#) in conjunction with merging logic), creating unique identifier variables or merging flags ensures the integrity of the integrated [dataset](#). This ability to integrate calculated and logical definitions into the DATA step is what makes [SAS](#) such a powerful tool for comprehensive data preparation.

The key takeaway for advanced users is that any valid [SAS](#) expression can be used on the right side of the assignment operator (=) to define a new variable. This includes intricate calculations, string manipulations using functions like `SUBSTR` or `CATX`, and date/time conversions. Always ensure that the resulting output type of the expression matches the intended variable type (e.g., a mathematical result should be assigned to a numeric variable), although [SAS](#) is generally forgiving and attempts automatic type conversion where possible, though explicit type management is always best practice.

Syntax Review and Best Practices for Efficient SAS Coding

To maximize efficiency and maintain code clarity, a few best practices should be followed when defining new [variables](#). First, adhere to clear and descriptive naming conventions. While [SAS](#) allows up to 32 characters for variable names, using short, meaningful names (like `avg_pts_rebs` instead of `x123`) significantly improves code readability and maintainability. When naming variables, avoid starting names with numbers and generally stick to letters, numbers, and underscores.

Second, always comment your code, especially before complex variable derivations. Clearly explaining the purpose of a new variable and the logic behind its calculation ensures that future analysts (or your future self) can quickly grasp the data transformation process. In the provided examples, we used `/*create dataset*/` and `/*view dataset*/`--this principle should be extended to every significant assignment statement.

Finally, remember the core syntax rules governing the two methods:

When creating variables from scratch, the [INPUT statement](#) must precede the [DATALINES statement](#), and the dollar sign (\$) is mandatory for [character variables](#).

When deriving variables from existing data, the [SET statement](#) must be used to load the existing [dataset](#) before any calculation statements.

Adherence to these guidelines ensures not only valid [SAS](#) code but also promotes a structured and standardized approach to data management that is essential in professional analytical environments.

Additional Resources

To further your skills in data preparation and manipulation, consider exploring the following advanced topics and related tutorials. Mastery of these concepts is essential for transitioning from basic data entry to complex statistical programming:

Understanding the difference between the [SET statement](#) and the [INPUT statement](#) in various contexts, including merging [datasets](#) and iterative processing.

Learning about explicit length and format statements (`LENGTH` and `FORMAT`) to precisely control how [variables](#) are stored and displayed.

Exploring the use of arrays for efficiently creating or transforming large groups of related [variables](#).

Reviewing documentation on temporary [variables](#) (those defined but not written to the output [dataset](#)) using the `DROP` or `KEEP` options within the DATA step.

The following tutorials explain how to perform other common tasks in [SAS](#):