

Use PROC SURVEYSELECT in SAS (With Examples)

Authored by
Mohammed looti

November 15, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Use PROC SURVEYSELECT in SAS (With Examples)*.
PSYCHOLOGICAL STATISTICS. Retrieved from
<https://statistics.arabpsychology.com/?p=1918>

Introduction: Harnessing [PROC SURVEYSELECT](#) for Precise Sampling in SAS

In the realm of statistical analysis, the validity of research findings hinges on obtaining a truly representative [sample](#) from a larger population. The powerful statistical software suite, [SAS](#), provides researchers with an indispensable procedure tailored specifically for this critical task: [PROC SURVEYSELECT](#). This procedure offers advanced capabilities for designing and executing various types of random sampling methodologies with exceptional precision and control, making it an essential tool for data professionals.

Generating accurate, unbiased [samples](#) is paramount for ensuring that any inferences drawn about the target population are statistically reliable and generalizable. Whether conducting complex observational studies, large-scale surveys, or controlled experiments, the integrity of the [sampling](#) design directly dictates the trustworthiness of the conclusions. [PROC SURVEYSELECT](#) provides a versatile framework capable of implementing everything from the most fundamental selection techniques to intricate, multi-stage [sampling](#) plans.

This comprehensive guide is designed to explore the practical implementation of [PROC SURVEYSELECT](#) through clear, runnable [SAS](#) examples. We will focus on three core [sampling](#) designs: [simple random sampling](#) (SRS), [stratified random sampling](#), and [clustered random sampling](#). By examining the underlying principles and corresponding code, analysts will gain the necessary skills to efficiently extract representative subsets of data for robust statistical analysis.

Foundational Principles of Statistical Sampling

Before utilizing the power of [SAS](#) code, it is essential to establish a firm grasp of the fundamental concepts underpinning the selection methods. Each sampling technique is designed to address specific data structures and research objectives, offering distinct advantages in terms of cost, effort, and precision.

[Simple Random Sampling \(SRS\)](#) stands as the baseline technique, ensuring that every possible subset of observations of a given size has an equal probability of being chosen. Critically, this means that every individual observation within the dataset has an equal chance of inclusion, resulting in an unbiased representation of the overall population. SRS is the preferred method when the population is considered relatively homogeneous and there are no critical subgroups that require explicit, proportionate representation.

In contrast, [Stratified Random Sampling](#) is employed when the population exhibits natural divisions into distinct subgroups, known as strata, which are internally uniform but highly variable when compared to one another. This method involves partitioning the population into these strata and then drawing a [simple random sample](#) independently from within each stratum. The primary

benefit of stratification is the guarantee that all important subgroups are adequately represented in the final sample, often leading to significantly more precise population estimates, especially when variability across the strata is high.

The third major design is [Clustered Random Sampling](#). This approach is highly practical when comprehensive lists of individual observations are unavailable or when the cost of sampling across a wide area is prohibitive. Here, the population is first divided into clusters (e.g., geographical areas or organizations), and then a random selection of these clusters is chosen. Typically, all observations contained within the selected clusters are included in the final sample. For example, selecting a sample of hospitals and surveying all nurses within those chosen hospitals.

A crucial best practice for any random [sampling](#) procedure is ensuring **reproducibility**. The [seed](#) option available in [PROC SURVEYSELECT](#) allows the user to specify a starting value for the random number generator. By fixing this [seed](#), running the code multiple times will yield the exact same random sample, which is fundamental for verification, peer review, and consistency in research findings.

Preparing the Hypothetical Dataset in SAS

To effectively demonstrate the capabilities of [PROC SURVEYSELECT](#), we will first create a sample dataset. This hypothetical dataset, named **my_data**, contains information pertaining to several basketball players. The dataset includes two key variables: **team** (identifying the player's team affiliation, which will serve as our stratification and clustering variable) and **points** (representing the individual player's scored points).

The following [SAS](#) code block illustrates the creation of this dataset using a standard **DATA** step. We populate the dataset efficiently using the **DATALINES** statement to input the player data. Following the data creation, the **PROC PRINT** procedure is executed. This step is vital for visually inspecting the dataset content and ensuring that the data structure is correct and ready for subsequent [sampling](#) operations.

```
/*create dataset for demonstration*/  
data my_data;  
input team $ points;  
datalines;  
A 12  
A 14  
A 22  
A 35  
A 40  
B 12
```

```
B 10
B 29
B 33
C 40
C 25
C 11
C 10
C 15
;
run;

/*view dataset structure and content*/
proc print data = my_data;
```

This essential preparatory step establishes the foundation for all three sampling methods we will explore. The output generated by **PROC PRINT** confirms that we have a dataset of 14 players distributed across three distinct teams (A, B, and C), providing the necessary heterogeneity to test the different sampling designs.

Obs	team	points
1	A	12
2	A	14
3	A	22
4	A	35
5	A	40
6	B	12
7	B	10
8	B	29
9	B	33
10	C	40
11	C	25
12	C	11
13	C	10
14	C	15

Implementing Simple Random Sampling (SRS) with PROC SURVEYSELECT

Our first demonstration focuses on [Simple Random Sampling \(SRS\)](#). As the most straightforward design, SRS is utilized when there is no need to account for subgroups within the population, and the objective is simply to select a subset where every observation has an equal probability of inclusion. This method serves as an excellent starting point for any sampling project.

To execute SRS using [PROC SURVEYSELECT](#), the core command is the specification of the selection method via the **METHOD=SRS** option. Furthermore, we use the **N=** option to explicitly define the required sample size, which in this case refers to the total number of individual observations to be selected. For rigorous scientific practice, we must include the **SEED=** option, setting a fixed value (here, 1) to ensure the exact random selection can be precisely reproduced at any time.

In this specific example, we instruct [SAS](#) to select 5 random observations from the **my_data** dataset. The resulting output dataset, named **my_sample**, will contain these five randomly chosen player records. The subsequent **PROC PRINT** procedure is then used to display the contents of **my_sample**, allowing for a verification of the randomly selected observations.

```
proc surveyselect data=my_data
out=my_sample
method=srs /*use simple random sampling*/
n=5 /*select a total of 5 observations*/
seed=1; /*set seed to make this example reproducible*/
run;

/*view sample*/
proc print data=my_sample;
```

The output image below visually confirms the result of the [simple random sample](#). Notice that the 5 selected observations were chosen without regard to their team affiliation, perfectly illustrating the principle of SRS where every row in the original dataset had an equal chance of being included.

Obs	team	points
1	B	10
2	B	33
3	C	40
4	C	25
5	C	11

Executing Stratified Random Sampling for Balanced Representation

When the population is heterogeneous, and specific groups must be guaranteed inclusion, [stratified random sampling](#) becomes the optimal methodology. This ensures that the final sample maintains balanced representation across critical subgroups, thereby minimizing sampling variance and improving the accuracy of population parameter estimates.

Using our basketball data, the **team** variable naturally defines our strata (Team A, B, and C). To implement this complex design in [PROC SURVEYSELECT](#), we introduce the indispensable **STRATA** statement, specifying the variable (**team**) by which the data should be partitioned. This tells [SAS](#) to treat these unique teams as independent sampling domains.

For this illustration, we aim to select an equal number of players from each team. We use the **N=2** option in conjunction with the **STRATA** statement, which instructs the procedure to select exactly 2 observations independently from *each* identified stratum. We retain **METHOD=SRS** because the selection process within each stratum is still a [simple random sample](#). Setting the [seed](#) ensures that this balanced selection is reproducible.

```
proc surveyselect data=my_data
out=my_sample
method=srs /*use simple random sampling within strata*/
n=2 /*select 2 observations from each strata*/
seed=1; /*set seed to make this example reproducible*/
strata team; /*specify variable to use for stratification*/
run;

/*view sample*/
proc print data=my_sample;
```

As clearly demonstrated in the accompanying image, the resulting sample contains six

observations: precisely 2 randomly chosen players from Team A, 2 from Team B, and 2 from Team C. This successful execution of [stratified random sampling](#) confirms balanced representation across all defined teams, achieving a sample that is highly representative of the underlying structure.

Obs	team	points	SelectionProb	SamplingWeight
1	A	35	0.4	2.5
2	A	40	0.4	2.5
3	B	29	0.5	2.0
4	B	33	0.5	2.0
5	C	25	0.4	2.5
6	C	10	0.4	2.5

Implementing Clustered Random Sampling for Efficiency

The final sampling method we will demonstrate is [clustered random sampling](#). This design is highly beneficial when individual listing of population units is impractical or when travel/logistical costs need to be minimized, as it groups observations together. Instead of sampling individuals, we sample entire groups or clusters, and then survey all members within those selected groups.

In the context of our basketball dataset, we define each team as a distinct cluster. Our goal is to randomly select a subset of these team clusters and include every player belonging to the chosen teams in our final sample. To achieve this within [PROC SURVEYSELECT](#), we employ the **CLUSTER** statement, specifying **team** as the variable that defines the clustering unit.

We use the **N=2** option, which, when paired with the **CLUSTER** statement, signifies the number of clusters to be selected, not the number of individual players. Therefore, we are instructing [SAS](#) to randomly select 2 out of the 3 available teams. The [seed](#) is fixed once more to ensure the chosen clusters are replicable. Upon running the procedure, **PROC PRINT** will list the full roster of players from the two randomly selected teams.

```
proc surveyselect data=my_data
out=my_sample
n=2 /*select a total of 2 clusters*/
seed=1; /*set seed to make this example reproducible*/
cluster team; /*specify variable to use for clustering*/
run;
```

```
/*view sample*/  
proc print data=my_sample;
```

The resulting sample, visible in the output image, includes every observation from Team A and Team B. These were the two clusters randomly chosen by the procedure. This output successfully demonstrates the mechanism of [clustered random sampling](#), where the final sample size is determined by the cumulative size of the selected groups.

Obs	team	points
1	A	12
2	A	14
3	A	22
4	A	35
5	A	40
6	B	12
7	B	10
8	B	29
9	B	33

Conclusion: Mastering Sampling Designs in SAS

The [PROC SURVEYSELECT](#) procedure within [SAS](#) is a versatile and fundamental component of the statistical workflow. Through the practical demonstrations provided--covering [simple random sampling](#), [stratified random sampling](#), and [clustered random sampling](#)--we have shown how this powerful tool facilitates the implementation of diverse random [sampling](#) designs. Correctly applying these methods ensures that researchers maintain the integrity and representativeness of their samples, leading directly to more robust and statistically defensible conclusions.

The choice of the appropriate [sampling](#) method is guided by several factors, including the specific research objectives, the inherent variability and structure of the population, and the logistical constraints of the study. [PROC SURVEYSELECT](#) offers the necessary flexibility to adapt to complex survey scenarios, guaranteeing that the data collection strategy is both efficient and statistically rigorous. As a consistent analytical best practice, always remember to fix the [seed](#) value to guarantee complete reproducibility of your random selection process.

For analysts seeking to utilize the full breadth of its capabilities, including probability proportional to

size (PPS) sampling and variance estimation techniques, the [official SAS documentation for PROC SURVEYSELECT](#) serves as the definitive resource for advanced options.

Additional Resources for SAS Proficiency

To further expand your proficiency in [SAS](#) programming and statistical tasks, explore these related tutorials:

How to Perform in SAS

Understanding in SAS

A Guide to with SAS