

Understanding Principal Component Analysis (PCA): A Step-by-Step Guide Using SAS

Authored by
Mohammed looti

November 14, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Understanding Principal Component Analysis (PCA): A Step-by-Step Guide Using SAS*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=1374>

The Core Principles of Principal Components Analysis (PCA)

[Principal Components Analysis \(PCA\)](#) is an indispensable and foundational statistical technique utilized extensively across modern [machine learning](#) and advanced statistical modeling workflows. The primary objective of PCA is not merely to simplify data, but to achieve rigorous [dimensionality reduction](#) of a complex [dataset](#) while judiciously preserving the maximal amount of actionable information held within the original structure. This powerful transformation process operates by converting a set of potentially correlated variables into a new, smaller set of entirely uncorrelated variables, which are formally known as **principal components**. These components are expertly constructed as [linear combinations](#) of the initial predictor variables, and they are strictly ordered so that the first few dimensions capture the largest possible degree of variation inherent in the underlying data structure.

The practical significance of PCA cannot be overstated in contemporary data analysis. It provides a vital mechanism for simplifying complex data interpretation and dramatically enhancing the computational efficiency and performance of subsequent analytical models. High-dimensional datasets frequently encounter challenges related to the "curse of dimensionality," rendering them difficult to visualize and often prohibitively expensive to process. By efficiently identifying the most significant underlying patterns and structures, PCA successfully mitigates these challenges, revealing hidden relationships that define the true structure of the data. Furthermore, it is supremely useful when analysts are faced with highly correlated variables, a situation that often leads to [multicollinearity](#) issues; the technique transforms these variables into a set of orthogonal (uncorrelated) components, thereby resolving these statistical dependencies.

The applications of PCA span an impressive breadth of domains, ranging from sophisticated image processing and facial recognition systems to intricate financial risk analysis and complex bioinformatics. For example, in the field of genomics, PCA can effectively cluster groups of genes that exhibit similar behavior across varied samples, providing critical insights into biological processes. Conversely, in market research, the technique offers the capability to condense responses from numerous survey questions into a handful of robust key factors that accurately represent core consumer preferences. Developing proficiency in both the execution and nuanced interpretation of PCA is therefore an absolutely essential skill set for any professional working with multivariate data, particularly within statistical programming environments like SAS.

Implementing PCA Using the SAS PROC PRINCOMP Statement

For professionals operating within the [SAS](#) statistical environment, executing Principal Components Analysis is a streamlined and efficient process facilitated by the dedicated [PROC PRINCOMP](#) statement. This highly specialized and powerful procedure is engineered to compute all core outputs of PCA, including the principal components, the critical [eigenvalues and](#)

[eigenvectors](#), while simultaneously offering a comprehensive suite of options for advanced analysis and customized data output creation. Its robust and flexible functionality establishes it as the definitive tool for [dimensionality reduction](#) within the SAS ecosystem, allowing analysts to quickly move from raw data to interpretable results.

The standard syntax required for utilizing `PROC PRINCOMP` is intuitive and straightforward, mirroring the structure of most SAS procedures. It requires users to specify the primary input dataset, meticulously select the specific variables designated for analysis, and define various options for the resulting output datasets. This flexibility is key, ensuring that the procedure can be precisely tailored to meet distinct analytical requirements, regardless of whether the user needs the transformed data, comprehensive statistical summaries, or specific visual representations needed for reporting. Below is the essential, foundational structure of the `PROC PRINCOMP` statement, illustrating how to initiate the analysis:

```
proc princomp data=my_data out=out_data outstat=stats;  
var var1 var2 var3;  
run;
```

A clear understanding of the purpose of each critical option within this syntax is essential for effective and customized implementation of PCA:

data: This parameter is absolutely crucial as it explicitly identifies the name of the input dataset that contains the numeric variables slated for analysis. Without a correctly and explicitly defined input dataset, `PROC PRINCOMP` cannot proceed with its underlying computations.

out: Employing this option triggers the creation of a new, supplementary dataset. This resultant output dataset will contain all the original variables derived from the input data, augmented by the newly computed principal component scores for every observation. These scores fundamentally represent the projection of each individual data point onto the newly defined principal component axes.

outstat: Utilizing this option generates a specialized secondary dataset that holds a substantial wealth of statistical information. This output includes essential descriptive statistics, the [correlation coefficients](#), and, most critically, the [eigenvalues and eigenvectors](#) derived directly from the analysis. This dataset is invaluable for achieving a comprehensive, detailed interpretation of the PCA results and for component selection.

var: Within the accompanying `VAR` statement, the user must list the specific numeric variables extracted from the input dataset that are intended for inclusion in the principal components analysis. It is imperative to carefully select only the variables that are relevant and meaningful to the defined PCA objective, as factor scaling can significantly impact the outcome.

The following comprehensive, step-by-step example is designed to practically illustrate the

application of the `PROC PRINCOMP` statement in a real-world scenario, clearly demonstrating its indispensable utility in successfully performing principal components analysis within the robust framework of [SAS](#).

Step-by-Step Execution: Data Preparation and Setup

Prior to initiating the core PCA, meticulous data preparation is an absolute prerequisite. This crucial initial phase involves either loading an existing dataset or constructing a new one that will be the subject of the analysis. For the purpose of this practical demonstration, we will construct a hypothetical dataset comprising performance statistics for 20 basketball players. This dataset will serve as the controlled foundation for our Principal Components Analysis, enabling us to efficiently identify the underlying, latent patterns in player performance metrics, such as potential groupings or specialization types.

This synthetic dataset incorporates three fundamental performance indicators: **points** scored, **assists** made, and **rebounds** grabbed. By incorporating and analyzing these three variables, our goal is to discern how they collectively contribute to the overall player profiles and to determine whether certain players exhibit comparable performance characteristics (e.g., scoring specialists versus rebound specialists). The construction of this dataset is accomplished efficiently using a standard SAS `DATA` step, immediately followed by a `PROC PRINT` statement, which is used to display and verify the generated data structure and ensure input accuracy.

```
/*create dataset*/  
data my_data;  
input points assists rebounds;  
datalines;  
22 8 4  
29 7 3  
10 4 12  
5 5 15  
35 6 2  
8 3 10  
10 4 8  
8 4 3  
2 5 17  
4 5 19  
9 9 4  
7 6 4  
31 5 3  
4 6 13
```

```
5 7 8  
8 8 4  
10 4 8  
20 4 6  
25 8 8  
18 8 3  
;  
run;
```

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

Once the SAS code responsible for the dataset creation has been successfully executed, the subsequent `PROC PRINT` statement will generate a tabular display of the contents of the `my_data` dataset. This crucial step permits a rapid visual inspection, confirming that the data has been accurately input, correctly structured, and is ready for the subsequent statistical analysis. This verification is a vital quality control measure before proceeding with complex multivariate procedures.

Obs	points	assists	rebounds
1	22	8	4
2	29	7	3
3	10	4	12
4	5	5	15
5	35	6	2
6	8	3	10
7	10	4	8
8	8	4	3
9	2	5	17
10	4	5	19
11	9	9	4
12	7	6	4
13	31	5	3
14	4	6	13
15	5	7	8
16	8	8	4
17	10	4	8
18	20	4	6
19	25	8	8
20	18	8	3

Performing and Reviewing the PCA Output

With the dataset successfully structured and prepared, the essential next step involves performing the principal components analysis using the powerful `PROC PRINCOMP` statement. This action initiates the transformation of our original, potentially correlated player performance variables into a new, optimized set of principal components, which are mathematically constructed to capture the maximum possible variance. The primary analytical goal here is to effectively distill the information contained within the variables **points**, **assists**, and **rebounds** into a smaller number of dimensions that are more statistically meaningful and interpretable.

To execute this analysis, we invoke `PROC PRINCOMP` and precisely specify `my_data` as the input source. We strategically utilize both the `OUT` and `OUTSTAT` options to generate two distinct and highly important output datasets: the `out_data` dataset will contain the original data merged with the new principal component scores, while the `stats` dataset will reliably house the comprehensive statistical summaries, including the critical [eigenvalues and eigenvectors](#)

necessary for interpreting the PCA results. The `VAR` statement explicitly selects the three basketball performance metrics for inclusion in the analysis, ensuring only relevant numeric variables are processed.

```
/*perform principal components analysis*/  
proc princomp data=my_data out=out_data outstat=stats;  
var points assists rebounds;  
run;
```

Upon the successful execution of this SAS code, `PROC PRINCOMP` generates a comprehensive output package. This output systematically includes various crucial elements: descriptive statistics, a detailed [correlation matrix](#) showing variable relationships, and the foundational eigenvalues and eigenvectors that constitute the mathematical core of the PCA. This extensive output is absolutely foundational for understanding the underlying structure of our data and assessing the proportional contribution of each derived principal component to the overall variance.

The initial segment of the output provides a detailed overview, commencing with descriptive statistics such as the mean and standard deviation for every variable included in the input. This is immediately followed by the [correlation matrix](#), which clearly reveals the pairwise linear relationships existing between our three performance variables--a high correlation here would validate the need for PCA. Finally, the eigenvalues and eigenvectors are presented, offering profound insights into the magnitude of variance explained by each component and elucidating the specific weighting (or component loading) of the original variables within the composition of those components. The eigenvectors are essential as they define the direction and orientation of the new axes in the reduced space.

The PRINCOMP Procedure

Observations	20
Variables	3

Simple Statistics			
	points	assists	rebounds
Mean	13.50000000	5.800000000	7.700000000
StD	10.06034424	1.765159900	5.120649630

Correlation Matrix			
	points	assists	rebounds
points	1.0000	0.2341	-.6232
assists	0.2341	1.0000	-.3855
rebounds	-.6232	-.3855	1.0000

Eigenvalues of the Correlation Matrix				
	Eigenvalue	Difference	Proportion	Cumulative
1	1.85089366	1.05544499	0.6170	0.6170
2	0.79544867	0.44179099	0.2651	0.8821
3	0.35365768		0.1179	1.0000

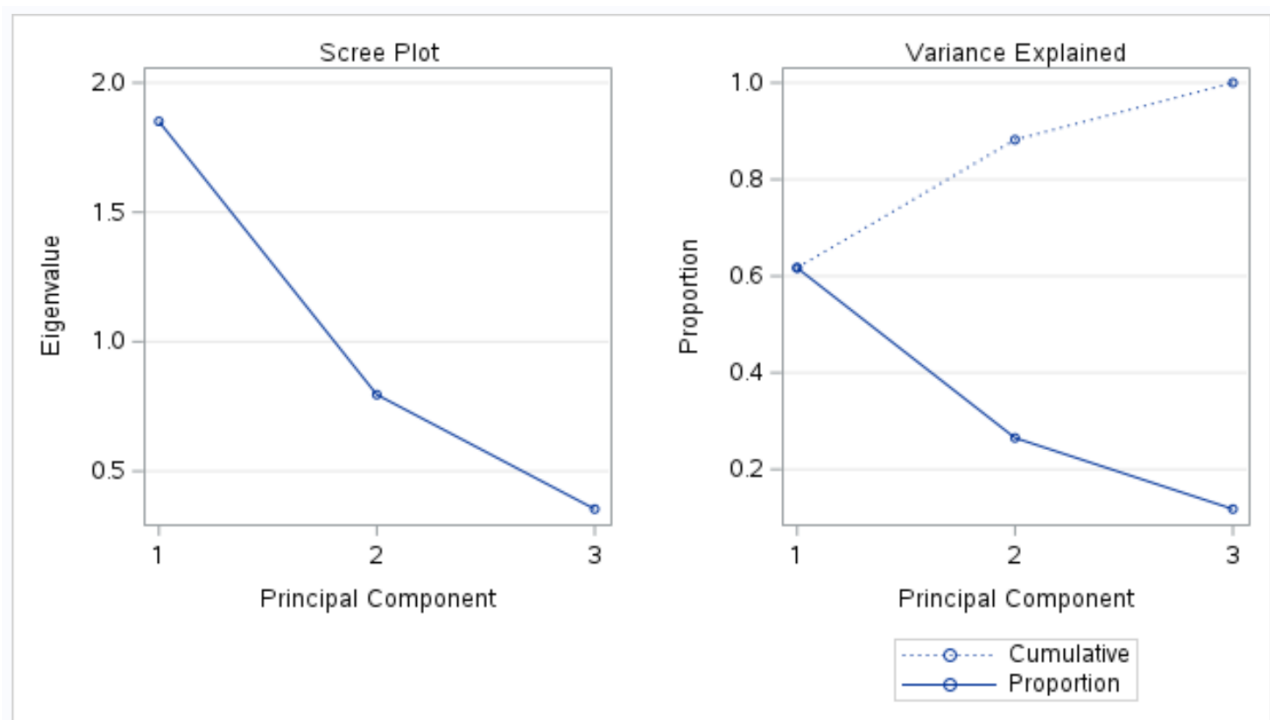
Eigenvectors			
	Prin1	Prin2	Prin3
points	0.603439	-.478471	0.637908
assists	0.460853	0.862106	0.210683
rebounds	-.650750	0.166848	0.740733

Interpreting Variance Explained and Component Selection

Interpreting the detailed output furnished by PROC PRINCOMP represents a vital stage in extracting meaningful and actionable insights from the multivariate data. The procedure delivers several key tables and diagnostic plots that are instrumental in fully understanding the principal components, evaluating their statistical significance, and quantifying the proportion of variance that each component successfully explains. This stage is crucial for deciding how many components should be retained for subsequent modeling or interpretation.

Immediately following the initial descriptive and statistical summaries, the output conventionally

includes essential graphical representations that are crucial for sound PCA interpretation. Specifically, the next portion of the output prominently displays a **Scree Plot** and a **Variance Explained** plot. These visualizations are fundamentally instrumental in assisting the analyst in determining the optimal, parsimonious number of principal components that should be retained for any subsequent, focused analysis, often guided by the Kaiser criterion (retaining components with eigenvalues greater than one) or by finding the "elbow" point in the Scree Plot.



The tabular output titled **Eigenvalues of the Correlation Matrix** is exceptionally informative. This table precisely quantifies the specific percentage of total variation that is accounted for by each successive principal component. Since each eigenvalue corresponds directly to a principal component, its numerical magnitude is a direct and powerful indicator of the amount of variance that specific component explains relative to the total variance in the dataset. Analyzing this table provides the quantitative foundation for component selection:

The first principal component (Prin1) clearly demonstrates its statistical dominance by explaining a substantial **61.7%** of the total variation observed in the dataset. This strongly indicates that a majority of the data's overall variability can be efficiently captured by this single underlying dimension, which likely represents overall player ability.

The second principal component (Prin2) contributes an additional significant portion, accounting for **26.51%** of the explained variance, further enriching our comprehension of the data's inherent structure.

The third principal component (Prin3) accounts for the remaining **11.79%** of the total variation

present in the data, offering only marginal explanatory power compared to the first two.

It is important to confirm that when these percentages are aggregated (61.7% + 26.51% + 11.79%), the sum precisely equals 100%. This provides confirmation that all the inherent variability present in the original dataset is fully explained by these three mathematically derived principal components--a predictable outcome given that we initiated the analysis with three original input variables. A common objective in PCA is to identify a minimal subset of components that collectively explain a sufficiently high cumulative percentage of the total variance (e.g., typically targeting 80% or 90%). In this specific case, the combination of the first two components alone explains 61.7% + 26.51%, totaling **88.21%** of the total variation. This result strongly suggests that these two components successfully capture the vast majority of the relevant information, making the third component potentially redundant and unnecessary for most practical applications of [dimensionality reduction](#).

Generating and Analyzing the PCA Biplot

To achieve a more profound, intuitive understanding of the complex relationships between the observations (players) and the original variables (metrics) within the newly established reduced dimensional space, it is highly recommended to construct a [biplot](#). The biplot serves as an exceptionally powerful visualization instrument that simultaneously maps every observation in the dataset and the original variables onto a single plane, typically formed by the first two principal components (Prin1 and Prin2). This visualization allows for a direct, simultaneous visual interpretation of how observations are naturally grouped and how the original variables contribute to the formation of these groupings, providing immediate insight into the component structure.

The biplot is essential for identifying patterns, such as clear clusters of similar observations and the precise directional influence exerted by the original variables on these clusters. For instance, variables that exhibit strong positive correlation will visually point in similar directions on the plot, while variables that are inversely correlated will point in diametrically opposite directions, clarifying the interpretation of the components. Furthermore, the length of the variable vectors provides an immediate visual indication of the strength and importance of their contribution to the corresponding principal components, linking the visual output back to the component loadings found in the [eigenvectors](#) table.

To successfully generate a biplot in [SAS](#), we must first ensure that our output data includes a unique identifier for each observation. This is accomplished by creating a new dataset, named `biplot_data`, which incorporates a column titled `obs` that numerically represents the row numbers of the original data. Following this preparatory step, the `PROC SGPLOT` statement is utilized. It leverages the calculated values from the first two principal components (`Prin1` and `Prin2`) to construct the scatter plot, with data points labeled accurately using our `obs` identifier, providing the

visual backbone of the biplot.

```
/*create dataset with column called obs to represent row numbers of original data*/
```

```
data biplot_data;
```

```
set out_data;
```

```
obs=_n_;
```

```
run;
```

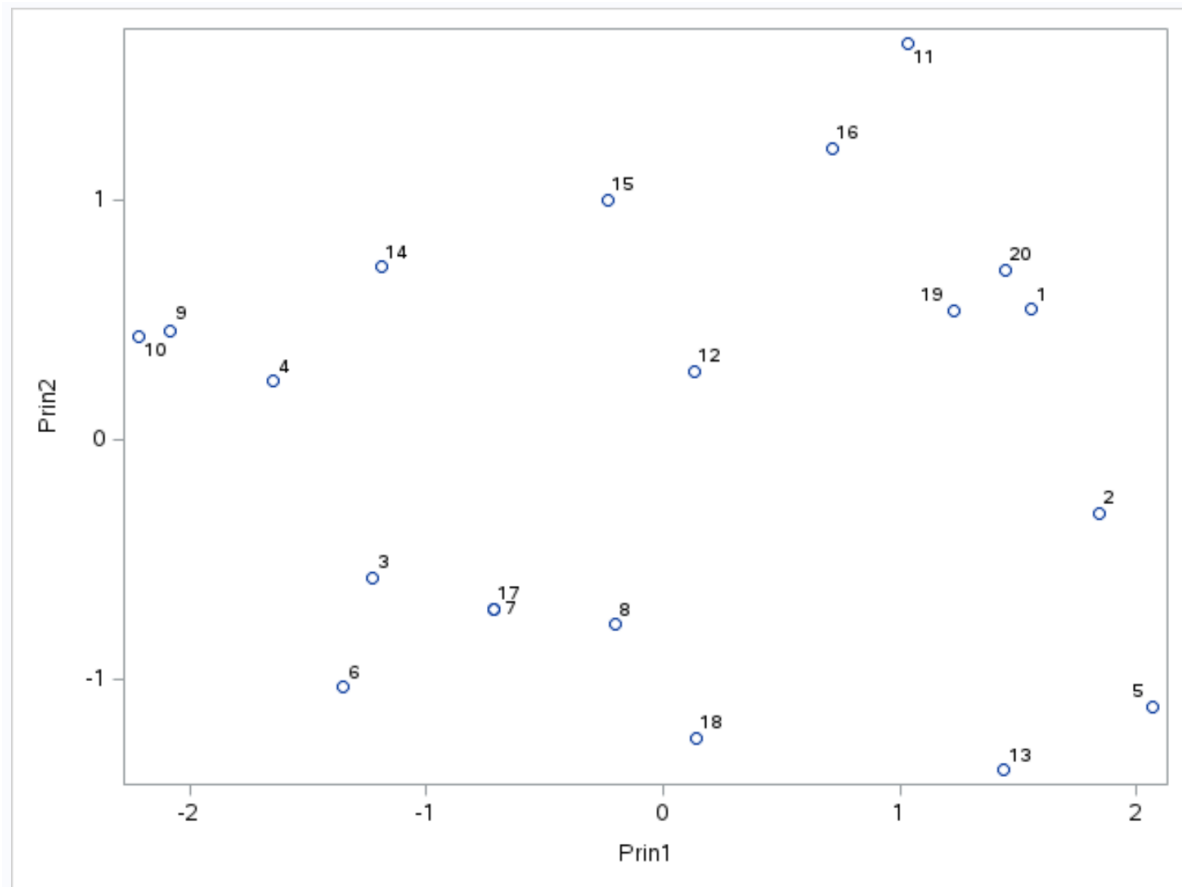
```
/*create biplot using values from first two principal components*/
```

```
proc sgplot data=biplot_data;
```

```
scatter x=Prin1 y=Prin2 / datalabel=obs;
```

```
run;
```

The resulting [biplot](#) offers a crucial visual representation of the intricate relationships embedded within our dataset. The x-axis corresponds to the dimension defined by the first principal component (Prin1), while the y-axis represents the dimension defined by the second principal component (Prin2). Each individual observation derived from our dataset is depicted as a discrete point within the plot, accurately labeled with its corresponding observation number, enabling analysts to trace specific data points back to their original records.



A critical element in interpreting the biplot is understanding the principle that observations positioned in close proximity to one another on the plot necessarily share similar characteristics across the original variables. For example, if we meticulously examine observations **#9** and **#10**, which are tightly clustered on the far left side of the plot, their proximity strongly suggests that their underlying values for **points**, **assists**, and **rebounds** are remarkably similar. We can validate this visual observation by cross-referencing the values from our original input dataset:

Observation #9: 2 points, 5 assists, 17 rebounds

Observation #10: 4 points, 5 assists, 19 rebounds

Indeed, the performance metrics for these two players are extremely close across all three indicators, which rigorously validates their tight proximity on the biplot. This confirms the biplot's efficacy as a tool for revealing inherent similarities and quantifiable differences among individual data points. Furthermore, recalling the analysis from the **Eigenvalues of the Correlation Matrix** table, we established that the first two principal components collectively account for **88.21%** of the total variation within the dataset. Given this high percentage, the biplot formed by these two components provides a highly reliable, accurate, and representative visualization, enabling analysts to confidently analyze the nuanced relationships and groupings among all observations.

Conclusion: Practical Utility and Methodological Caveats

Principal Components Analysis remains an indispensable statistical methodology for effectively simplifying large, complex datasets and extracting profound, actionable insights. As meticulously demonstrated through the practical application of the `PROC PRINCOMP` statement in SAS, PCA provides a systematic, rigorous approach to [dimensionality reduction](#), making the data significantly more manageable for both visualization and subsequent advanced modeling. By successfully transforming highly correlated variables into a concise set of uncorrelated **principal components**, PCA aids significantly in identifying latent patterns and underlying structures that might otherwise remain entirely obscured within high-dimensional data spaces.

The end-to-end process--spanning from initial data preparation and the execution of `PROC PRINCOMP`, through the critical interpretation of [eigenvalues and eigenvectors](#), [scree plots](#), variance explained plots, and biplots--establishes a comprehensive and reliable workflow for fully understanding data variability. The inherent ability of PCA to condense crucial information, as clearly seen in our basketball player example where only two components captured over 88% of the total variance, underscores its exceptional practical utility across diverse analytical contexts. This successful dimensionality reduction not only facilitates far clearer interpretation but also often leads to significant improvements in the stability and efficiency of predictive models by reducing noise and redundancy.

While PCA is a powerful tool, it is essential to acknowledge its primary methodological caveats. Firstly, the resulting principal components are inherently [linear combinations](#) of the original variables, which can sometimes render their direct, practical interpretation less intuitive compared to analyzing the original features. Secondly, PCA fundamentally assumes linearity in the data structure and is notably sensitive to the scaling of variables. Therefore, standardizing the data (a feature easily accessible using the `STAND` option within `PROC PRINCOMP`) is frequently a highly recommended preprocessing step. This standardization ensures that variables possessing naturally larger variances do not disproportionately dominate or influence the composition of the derived components, leading to a more equitable and robust analysis. Mastering PCA within the SAS environment equips analysts with a robust and versatile method for rigorous exploratory data analysis and effective feature engineering.

Additional Resources and Further Reading

To further refine your proficiency in SAS and explore other complex data analysis tasks, we recommend consulting these authoritative tutorials and documentation sources:

SAS Official Documentation for PROC PRINCOMP: [Link](#)

Wikipedia on Principal Components Analysis: [Link](#)