

# Calculate R-Squared in SAS

Authored by  
**Mohammed looti**

November 15, 2025

## RECOMMENDED CITATION

Mohammed looti (2025). *Calculate R-Squared in SAS*. PSYCHOLOGICAL STATISTICS.  
Retrieved from <https://statistics.arabpsychology.com/?p=1959>

## The Crucial Role of R-Squared in Model Assessment

In the field of [statistical analysis](#), and particularly when building predictive models through [regression analysis](#), understanding model fit is paramount. The metric known as **R-squared** ( $R^2$ ), or the Coefficient of Determination, stands out as a fundamental measure for assessing how well a chosen [statistical model](#) aligns with the observed data. It provides an immediate, intuitive summary of the relationship between your input variables and the outcome you are trying to predict.

R-squared quantifies the proportion of the total variation in the response variable that is successfully explained by the independent predictor variables included in the model. Essentially, it answers the question: how much of the scatter or [variance](#) in the outcome can we account for using the inputs? A high R-squared value indicates that the regression model captures a significant portion of the variability in the dependent variable, suggesting a strong explanatory relationship between the predictors and the outcome.

Since **R-squared** is a proportion, its value is strictly bounded, ranging from 0 to 1 (or 0% to 100%). Interpreting these boundaries is vital for model diagnostics:

A value of **0** signifies that the predictor variable(s) offer no explanatory power whatsoever; the model performs no better than simply using the mean of the response variable to predict future outcomes. In practical terms, this means the chosen predictors provide zero insight into the outcome's variability.

A value of **1** (or 100%) represents a perfect fit. In this highly idealized scenario, all observed data points lie exactly on the regression line, meaning the model perfectly explains all the [variance](#) in the response variable with zero unexplained error.

While an exceptionally high R-squared is often desirable, it should never be interpreted in isolation. It is a measure of fit, not necessarily a guarantee of accuracy or causality. Over-reliance on a high R-squared can sometimes indicate [overfitting](#), where the model explains the sample data perfectly but fails to generalize to new, unseen data. Therefore, context, domain expertise, and a suite of additional diagnostic tests are always required for thorough model validation.

## A Practical Case Study: Calculating R-Squared using SAS

To solidify the theoretical understanding of R-squared, we will transition to a practical demonstration using [SAS](#) (Statistical Analysis System), a leading software package for statistical computing and data management. This step-by-step guide walks you through the entire process: from data creation to model fitting and, finally, the precise extraction of the R-squared metric for a [simple linear regression model](#).

Our objective for this exercise is to model the relationship between the time students dedicate to studying and their resulting final exam scores. This is a classic application of regression where we hypothesize a linear connection. We will treat the total hours studied as our primary [predictor variable](#) and the final exam score as the [response variable](#).

The following procedures illustrate the efficiency and clarity of using SAS procedures to generate regression output and isolate key performance indicators like **R-squared**. By executing the code examples provided, you will gain hands-on experience in generating and interpreting this essential metric within a professional statistical environment, ensuring accurate and reproducible results.

## Step 1: Defining and Populating the Dataset in SAS

The initial step in any SAS analysis involves defining the [dataset](#) that will be used for modeling. We will create a small sample dataset containing 15 hypothetical records to simulate student performance. These records consist of two essential variables: `hours`, representing the total study time, and `score`, representing the final exam grade.

The subsequent SAS code block demonstrates the fundamental syntax for data entry. The `DATA` step names the new dataset `exam_data`. The `INPUT` statement defines the variable names and their order. The `DATALINES` statement signals to SAS that the data immediately follows in an inline format. After the data is entered, the `RUN` statement executes the data creation step. Furthermore, we employ `PROC PRINT` to visually confirm that the data has been loaded correctly before proceeding with the analysis. This ensures data integrity prior to modeling.

```
/*create dataset*/  
data exam_data;  
input hours score;  
datalines;  
1 64  
2 66  
4 76  
5 73  
5 74  
6 81  
6 83  
7 82  
8 80  
10 88  
11 84  
11 82
```

```
12 91
12 93
14 89
;
run;

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

Obs	hours	score
1	1	64
2	2	66
3	4	76
4	5	73
5	5	74
6	6	81
7	6	83
8	7	82
9	8	80
10	10	88
11	11	84
12	11	82
13	12	91
14	12	93
15	14	89

The image above confirms the successful execution of the data step, showing the `exam_data` dataset properly structured with the `hours` and `score` variables, verifying that the sample data is ready for the subsequent regression procedure.

## Step 2: Executing Simple Linear Regression using PROC REG

With the data prepared, the next crucial step is to fit the statistical model. In SAS, the standard and most powerful tool for performing [regression analysis](#) is the **PROC REG** procedure. This procedure is optimized for linear modeling and provides extensive diagnostic and output options, making it central to statistical analysis in the SAS environment.

The code required for running the regression is straightforward: the `PROC REG` statement specifies the procedure and the input data (`data=exam_data`). Crucially, the `MODEL` statement formally

defines the regression relationship: `score = hours`. This syntax clearly establishes `score` as the dependent (response) variable and `hours` as the independent (predictor) variable. The execution of this step generates a detailed report containing all standard regression output tables, including parameter estimates and the ANOVA summary.

```
/*fit simple linear regression model*/
```

```
proc reg data=exam_data;
```

```
model score = hours;
```

```
run;
```

The REG Procedure  
Model: MODEL1  
Dependent Variable: score

Number of Observations Read	15
Number of Observations Used	15

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	847.26698	847.26698	63.91	<.0001
Error	13	172.33302	13.25639		
Corrected Total	14	1019.60000			

Root MSE	3.64093	R-Square	0.8310
Dependent Mean	80.40000	Adj R-Sq	0.8180
Coeff Var	4.52852		

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	65.33395	2.10599	31.02	<.0001
hours	1	1.98237	0.24796	7.99	<.0001

Reviewing the comprehensive output generated by [PROC REG](#), we can quickly locate the dedicated section for model summary statistics. The raw **R-squared** value is prominently presented as **0.8310**. This result indicates that approximately 83.10% of the observed [variance](#) in the students' final exam scores is statistically explained by the number of hours they dedicated to studying, demonstrating a strong explanatory power for this specific [model](#).

### Step 3: Streamlining Output to Extract Only the R-Squared Metric

In advanced analytical environments or when developing automated reports, generating the full, multi-table output from **PROC REG** can be inefficient. Often, analysts require only a single metric, such as **R-squared**, to feed into subsequent calculations or summary reports. SAS provides powerful options to suppress comprehensive output and store specific statistics directly into a new dataset for focused processing.

To achieve this focused extraction, we introduce two critical options within the `PROC REG` statement: `NOPRINT` and `OUTEST=`. The `NOPRINT` option is essential as it completely suppresses the standard output tables generated by the procedure, ensuring a cleaner execution log. The `OUTEST=outest` option instructs SAS to create a new output [dataset](#) (named `outest` in this example) specifically designed to hold estimated parameters and model fit statistics.

Furthermore, we modify the `MODEL` statement by adding the `/ RSQUARE` option. This explicitly tells the procedure to calculate and store the R-squared value within the `outest` dataset under the standardized variable name `_RSQ_`. After running the regression procedure, the `QUIT` statement properly terminates the interactive procedure. A final `PROC PRINT` step, utilizing `var _RSQ_`, is then used to display only the desired R-squared value stored in the new dataset, isolating the required metric efficiently.

```
/*fit simple linear regression model*/  
proc reg data=exam_data outest=outest noprint;  
model score = hours / rsquare;  
run;  
quit;
```

```
/*print R-squared value of model*/  
proc print data=outest;  
var _RSQ_;  
run;
```

Obs	_RSQ_
1	0.83098

The resulting output confirms the successful extraction, yielding the precise **R-squared** value of

**0.83098.** This methodology is indispensable when integrating SAS statistical results into automated workflows, ensuring that only necessary metrics are generated and processed, thereby optimizing computational efficiency.

## Conclusion: Beyond R-Squared in Comprehensive Model Evaluation

Whether you prefer the complete regression report or the streamlined extraction method, calculating **R-squared** in [SAS](#) remains a fundamental skill for data professionals. As demonstrated, R-squared provides a clear, digestible measure of the explanatory power of your regression equation, linking the proportion of the [variance](#) in the outcome variable to the input predictors.

However, statistical proficiency demands looking beyond this single metric. A comprehensive [regression analysis](#) requires evaluating the model from multiple angles. For instance, the [adjusted R-squared](#) is often preferred because it penalizes the inclusion of unnecessary predictor variables, providing a more honest assessment of fit, especially in multiple regression scenarios.

Furthermore, examining the statistical significance of individual coefficients through [p-values](#) and assessing prediction error using metrics like [RMSE](#) (Root Mean Squared Error) is crucial for a complete understanding of model reliability. By integrating the calculation of **R-squared** alongside these other diagnostic tools, you ensure that your statistical models are not only a good fit for the historical data but are also robust, parsimonious, and reliable for future predictions. Mastery of these techniques in SAS allows for the production of high-quality, trustworthy analytical results.

## Additional Resources for Advancing Your SAS Skills

To further enhance your command of statistical programming and data management within SAS, we recommend exploring tutorials that delve deeper into advanced procedures and data manipulation techniques. Continuous learning in these areas is key to leveraging the full capabilities of the SAS platform for complex data science problems.