

Understanding Sum of Squares: Calculating SST, SSR, and SSE in R for Regression Analysis

Authored by
Mohammed loot

November 5, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Understanding Sum of Squares: Calculating SST, SSR, and SSE in R for Regression Analysis*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=10835>

When assessing the explanatory power and overall suitability of a [statistical model](#), particularly within the domain of [linear regression](#), analysts must rely on precise mathematical measures that quantify the [variance](#) inherent in the observed data. These fundamental statistical metrics are essential tools, enabling us to rigorously determine the extent to which the total variability observed in the response variable can be successfully accounted for by the predictive variables included in the model.

The concept of Sums of Squares is foundational to this evaluation process. They establish a robust framework for partitioning the total variation of the dependent variable into two crucial and distinct components: the variation explained by the proposed regression model (signal) and the residual, unexplained variation (noise or error). To measure how effectively a regression model aligns with and predicts a given dataset, we consistently utilize three critical values:

Sum of Squares Total (SST)

Sum of Squares Regression (SSR)

Sum of Squares Error (SSE)

Understanding Variance Partitioning in Regression

The framework established by the Sums of Squares is indispensable in statistical inference because it provides the necessary components for calculating the crucial [F-statistic](#) and, most notably, the Coefficient of Determination, also known as [R-squared](#). Taken together, these metrics provide a comprehensive narrative regarding the model's predictive strength and its overall explanatory capacity within the context of the data.

The relationship governing these three measures is an identity defined by the equation: **SST = SSR + SSE**. This straightforward mathematical relationship clearly signifies that the total inherent variability (SST) within the response variable must equal the sum of the variability successfully captured by the regression line (SSR) and the variability that stubbornly remains unexplained by the model (SSE). This partitioning provides the mathematical basis for analyzing model fit.

Achieving an accurate interpretation of the regression model hinges upon a precise understanding of how each of these components is calculated. The following section formally defines these metrics, focusing on their specific mathematical roles in the process of variance partitioning.

Deconstructing the Core Metrics: SST, SSR, and SSE

Each Sum of Squares term serves to quantify the distance between specific data points and either the mean value or the predicted values from the model. Crucially, these differences are squared. This methodology serves two primary purposes: first, it ensures that all calculated measures of deviation are positive, and second, it heavily penalizes larger deviations, leading to a more robust

and sensitive measure of overall variability.

1. Sum of Squares Total (SST) - This metric represents the total, raw [variation](#) present in the dependent response variable, Y. It is formally calculated by summing the squared differences between each individual observation (y_i) and the overall grand mean of the response variable (\bar{y}). The SST establishes the absolute baseline measure of variation that the subsequent model attempts to explain.

$$SST = \sum (y_i - \bar{y})^2$$

2. Sum of Squares Regression (SSR) - Often referred to as the Sum of Squares Explained, SSR quantifies the portion of the response variable's variation that is successfully captured and explained by the fitted regression model. Mathematically, it is the sum of squared differences between the predicted values derived from the model (\hat{y}_i) and the mean of the response variable (\bar{y}). A high SSR, when evaluated relative to the SST, is the key indicator of a stronger model fit.

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

3. Sum of Squares Error (SSE) - Also known as the Residual Sum of Squares, SSE measures the unexplained [variation](#) within the dataset. It represents the specific amount of variability in the dependent variable that the chosen predictor variable(s) failed to account for. This metric is calculated as the sum of squared differences between the observed data points (y_i) and the corresponding predicted data points (\hat{y}_i). These differences are fundamentally the model's [residuals](#).

$$SSE = \sum (y_i - \hat{y}_i)^2$$

The following detailed, step-by-step example transitions these theoretical concepts into practical application, demonstrating precisely how to calculate each of these metrics for a specific [linear regression](#) model using the powerful [R programming language](#).

Step 1: Setting Up the Environment and Preparing the Data

Before any statistical model can be rigorously fitted and analyzed, it is paramount to establish a clear and structured dataset that accurately reflects the relationship targeted for study. For the purpose of this practical demonstration, we will analyze the quantifiable relationship between the hours students dedicate to studying and the subsequent examination scores they achieve. This scenario represents a textbook application of simple [linear regression](#), where "hours studied" serves as the independent variable (predictor) and "exam score" is the dependent variable (response).

We initiate the process by constructing a data frame within R. This structure contains simulated

data representing the study hours and corresponding exam scores for a cohort of 20 distinct academic students. This preparatory step ensures that we have a clean, structured environment ready for model training and the subsequent explicit calculation of the Sums of Squares metrics.

Create the data frame containing hours studied and exam scores

```
df <- data.frame(hours=c(1, 1, 1, 2, 2, 2, 2, 2, 3, 3,
3, 4, 4, 4, 5, 5, 6, 7, 7, 8),
score=c(68, 76, 74, 80, 76, 78, 81, 84, 86, 83,
88, 85, 89, 94, 93, 94, 96, 89, 92, 97))
```

```
# Display the first six observations of the data frame to confirm data integrity
head(df)
```

```
hours score
1 1 68
2 1 76
3 1 74
4 2 80
5 2 76
6 2 78
```

A preliminary review of the generated dataset confirms the expected positive correlation: generally, as the total number of hours dedicated to studying increases, the corresponding exam scores also exhibit an upward trend. Our subsequent objective is to precisely quantify this relationship and assess its predictive strength using a formal regression model.

Step 2: Fitting the Linear Regression Model

The next essential phase in the analysis requires fitting the prepared data to a simple linear regression model. In the [R programming language](#), this is accomplished efficiently using the highly versatile `lm()` function, which is the standard command for creating "linear models." We must clearly specify the model formula, designating `score` as the response variable and `hours` as the solitary predictor variable.

By executing this model fitting process, we instruct R to determine the unique "best-fit" line. This line is mathematically optimized to minimize the residual Sum of Squares Error (SSE). The output establishes the precise intercept and slope coefficients that mathematically define the relationship between studying time and academic performance.

Fit the simple linear regression model

```
model <- lm(score ~ hours, data = df)
```

```
# View the detailed summary of the model fit
summary(model)
```

Call:

```
lm(formula = score ~ hours, data = df)
```

Residuals:

```
Min 1Q Median 3Q Max
```

```
-8.6970 -2.5156 -0.0737 3.1100 7.5495
```

Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
```

```
(Intercept) 73.4459 1.9147 38.360 < 2e-16 ***
```

```
hours 3.2512 0.4603 7.063 1.38e-06 ***
```

```
---
```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 4.289 on 18 degrees of freedom

Multiple R-squared: 0.7348, Adjusted R-squared: 0.7201

F-statistic: 49.88 on 1 and 18 DF, p-value: 1.378e-06

The resulting model summary provides several pieces of crucial statistical information. Key outputs include the estimated coefficients: an intercept of 73.4459 and a slope coefficient for hours of 3.2512. This slope indicates that for every additional hour a student studies, their exam score is expected to increase by approximately **3.25 points**. Furthermore, the summary explicitly provides the [R-squared](#) value, a metric we will now manually verify using our explicit Sums of Squares calculations.

Step 3: Explicit Calculation of SST, SSR, and SSE

While the standard R model summary conveniently provides the [R-squared](#) value, performing the explicit calculation of SST, SSR, and SSE is invaluable for fully grasping the mechanism of variance separation and ensuring methodological transparency. We utilize R's powerful built-in functions for this process, specifically `fitted()`, which accurately retrieves the predicted values (?i), alongside standard aggregation functions like `sum()` and `mean()`.

The sequence of calculation typically prioritizes SSE first, as it relies on the difference between the actual observed scores and the model's predicted scores. Subsequently, we calculate SSR using the difference between the predicted scores and the overall mean score. Finally, SST is derived by simply summing SSR and SSE, a step that simultaneously serves as a critical confirmation of the fundamental variance partitioning identity ($SST = SSR + SSE$).

Calculate Sum of Squares Error (SSE): Sum of squared differences between fitted values and actual scores

```
sse <- sum((fitted(model) - df$score)^2)
```

```
sse
```

```
331.0749
```

Calculate Sum of Squares Regression (SSR): Sum of squared differences between fitted values and the mean score

```
ssr <- sum((fitted(model) - mean(df$score))^2)
```

```
ssr
```

```
917.4751
```

Calculate Sum of Squares Total (SST): The sum of SSR and SSE

```
sst <- ssr + sse
```

```
sst
```

```
1248.55
```

The output of these rigorous calculations provides the essential quantitative metrics required for a thorough assessment of our model's performance. The final calculated values are summarized below:

Sum of Squares Total (SST): 1248.55

Sum of Squares Regression (SSR): 917.4751

Sum of Squares Error (SSE): 331.0749

To ensure absolute correctness, we perform a crucial verification step, confirming that the fundamental mathematical identity holds true:

$$SST = SSR + SSE$$
$$1248.55 \approx 917.4751 + 331.0749$$

This successful confirmation validates the precision of our calculations, robustly demonstrating that the total [variance](#) in the dataset has been correctly and entirely partitioned into its explained (SSR) and unexplained (SSE) components.

Interpreting R-squared and Model Fit

The paramount practical application of calculating SST and SSR is deriving the Coefficient of Determination, universally recognized as R-squared. This highly critical statistic measures the

proportion of the total variance in the dependent variable that can be reliably predicted or accounted for by the independent variable(s) included in the model.

The mathematical definition of [R-squared](#) is the simple ratio of the variance successfully explained by the model (SSR) to the total variance present in the data (SST):

$$R\text{-squared} = SSR / SST$$

Applying the precise values calculated from our dedicated R analysis in Step 3, we manually compute the R-squared value:

$$R\text{-squared} = 917.4751 / 1248.55$$

$$R\text{-squared} \approx 0.7348$$

This result achieves a perfect match with the "Multiple R-squared" value provided in the comprehensive summary output from Step 2 (0.7348). This high value strongly suggests that the model provides a highly effective fit for the relationship between study hours and exam scores.

Specifically, this interpretation leads to the conclusion that **73.48%** of the total variation observed in the students' exam scores can be effectively and directly explained by the variation in the number of hours they dedicated to studying. The remaining 26.52% of the variation is attributed either to random error (SSE) or to unmeasured confounding factors not incorporated into this specific [statistical model](#), such as differences in prior academic knowledge, levels of test anxiety, or variations in teaching quality.

Enhancing Regression Diagnostics

While the explicit manual calculation of SST, SSR, and SSE is crucial for mastering the underlying statistical principles, practitioners routinely utilize automated tools for rapid verification and for conducting more complex analyses involving multiple predictor variables. This detailed breakdown, however, offers a necessary transparent view into the foundational structure and performance assessment of the model.

For ensuring the integrity and validity of the statistical inferences drawn from any model, further exploration into advanced regression diagnostics is highly recommended. This includes critically analyzing the distribution characteristics of the [residuals](#) and systematically checking key assumptions like homoscedasticity. Utilizing the [R programming language](#) facilitates rapid iteration and testing across various model specifications and diagnostic checks.

Finally, for users seeking automated verification of these core metrics, numerous specialized statistical software packages and online calculators exist. These tools often automatically compute and integrate SST, SSR, and SSE directly within the primary model summary output, mirroring the

efficient functionality demonstrated by R's `summary(model)` function.