

# How to perform LOESS Regression in R (With Example)

Authored by  
**stats writer**

November 29, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to perform LOESS Regression in R (With Example)*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=101512>

LOESS regression is a sophisticated non-parametric regression technique designed to capture complex, non-linear relationships within data without requiring the specification of a global function structure. This method operates by fitting multiple, localized linear regression models across subsets of the data, weighting observations based on their proximity. Specifically, the technique uses locally weighted polynomial regression to estimate the smooth curve, making it exceptionally flexible for exploratory data analysis and smoothing time series data. In the R programming language, this is efficiently implemented via the `loess` function found within the base `stats` package, providing analysts with a powerful tool for visualization and modeling.

Unlike traditional parametric models which impose constraints like linearity (e.g., standard ordinary least squares), LOESS regression allows the relationship between the response variable and one or more predictor variables to change across the data space. This flexibility is achieved through a smoothing parameter, known as the span, which dictates the size of the local neighborhood used for each weighted fit. For instance, to fit a basic LOESS model to the built-in `mtcars` dataset in R, examining how miles per gallon (`mpg`) relates to vehicle weight (`wt`), one would use the straightforward syntax: `loess(mpg ~ wt, data = mtcars)`. Understanding the selection and optimization of the span parameter is crucial to generating a robust and meaningful fit, which is the focus of the subsequent steps.

**LOESS regression**, often referred to as local regression or locally estimated scatterplot smoothing, provides a powerful, data-driven approach to modeling smooth trends. By concentrating on local subsets of data, it effectively mitigates the rigid assumptions associated with global parametric models, offering a nuanced view of the underlying data structure. The following guide details a comprehensive, step-by-step methodology for executing and optimizing LOESS regression models within the R environment.

The primary challenge in applying this method lies in selecting the optimal smoothing parameter (span) to achieve a desirable balance between flexibility (low bias) and stability (low variance). Too tight a span leads to overfitting, capturing noise; too wide a span leads to underfitting, missing critical local trends. We will demonstrate how to manage this crucial tradeoff using statistical cross-validation techniques available in R's specialized packages.

## Understanding LOESS Regression and Its Mechanism

LOESS stands out due to its reliance on local fitting. When predicting a value for a specific data point, the method assigns weights to all neighboring points, prioritizing those closer to the point of interest. This weighting mechanism typically employs a tricube function, which ensures that influence decays rapidly as distance increases. A localized polynomial (usually first or second degree) is then fitted to this weighted neighborhood. This process is repeated across all points in the dataset, effectively "stitching together" many local models to form a single, smooth regression

curve.

The core philosophy of this approach is that while the relationship might be highly complex globally, it can be approximated well by a simple polynomial within a small, localized region. This adaptability allows LOESS regression to handle diverse data patterns, including sudden changes in slope or curvature, which would be challenging for a single global model to capture accurately. However, due to its computational intensity and reliance on the selection of the span, it is often best suited for datasets with moderate size.

The implementation in R is highly accessible, primarily leveraging the `loess()` function. This function abstracts the complex mathematical processes--including the selection of the polynomial degree and the iterative fitting algorithm--allowing the user to focus mainly on defining the relationship (via the formula syntax, e.g., `y ~ x`) and, most importantly, tuning the smoothing parameter.

## Practical Implementation: Defining the Dataset

To illustrate the process of fitting and optimizing a LOESS model, we must first establish a suitable dataset. For effective demonstration, we often choose data that exhibits a clear non-linear trend, making it unsuitable for standard linear modeling. The following steps involve creating a synthetic data frame in R, which captures a distinct curvilinear relationship between the predictor variable `x` and the response variable `y`.

This initial step ensures that we have a controlled environment to observe how different parameters affect the final fitted curve. Pay close attention to the structure of the data frame, as the subsequent fitting steps depend directly on these defined variables. Using a synthetic dataset allows us to isolate the effects of the smoothing span without the complexity of external confounding factors present in real-world data.

We begin by defining the data frame `df`, which contains 14 observations. This modest size is typical for initial LOESS modeling exercises, allowing us to easily visualize the effects of the smoothing process. Reviewing the first few rows confirms the successful creation and structure of our working data set before proceeding to the modeling phase.

### Step 1: Create the Data

First, we create the following data frame in R, which will serve as the foundation for our localized regression analysis:

```
#view DataFrame
```

```
df <- data.frame(x=c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14),
```

```
y=c(1, 4, 7, 13, 19, 24, 20, 15, 13, 11, 15, 18, 22, 27))
```

```
#view first six rows of data frame
```

```
head(df)
```

```
x y
```

```
1 1 1
```

```
2 2 4
```

```
3 3 7
```

```
4 4 13
```

```
5 5 19
```

```
6 6 24
```

This data frame clearly shows that as  $x$  increases,  $y$  initially rises sharply, then dips, and finally rises again, indicating a complex, non-linear pattern that demands a flexible modeling approach like LOESS. The success of the subsequent modeling relies heavily on how we define the parameters of the local fitting algorithm to accurately trace this underlying pattern without over-smoothing or overfitting.

## Step 2: Fit Several LOESS Regression Models

The core of applying LOESS involves utilizing the `loess()` function and critically examining the effect of the **span** parameter. The span defines the proportion of the data used to estimate the fit at each point. For example, a span of 0.5 means that 50% of the total data points nearest to the target point are included in that local regression calculation. Choosing the appropriate span is perhaps the most important decision in LOESS modeling, as it directly controls the smoothness of the fitted curve and dictates the trade-off between bias (smoothness) and variance (sensitivity to local noise).

We will fit three distinct LOESS models to our synthetic dataset, employing span values of 0.5, 0.75, and 0.9. By selecting a range of values--from a more localized fit (0.5) to a smoother, more global fit (0.9)--we can visually and quantitatively assess how the span influences the resulting prediction curve. A smaller span generally yields a highly flexible curve that closely follows the data points, which risks overfitting, especially if the data contains significant noise. Conversely, a larger span creates a much smoother curve, potentially masking important local features but offering greater stability.

After fitting the models, we use the `predict()` function to generate the smoothed values (`smooth50`, `smooth75`, `smooth90`). These predicted values represent the fitted regression curve for each respective span. Plotting these three curves simultaneously provides an immediate visual

comparison, making the effect of the span parameter intuitively clear to the analyst. This visualization is essential for developing an understanding of how hyperparameter selection impacts model performance.

**#fit several LOESS regression models to dataset**

```
loess50 <- loess(y ~ x, data=df, span=.5)
```

```
smooth50 <- predict(loess50)
```

```
loess75 <- loess(y ~ x, data=df, span=.75)
```

```
smooth75 <- predict(loess75)
```

```
loess90 <- loess(y ~ x, data=df, span=.9)
```

```
smooth90 <- predict(loess90)
```

**#create scatterplot with each regression line overlaid**

```
plot(df$x, df$y, pch=19, main='Loess Regression Models')
```

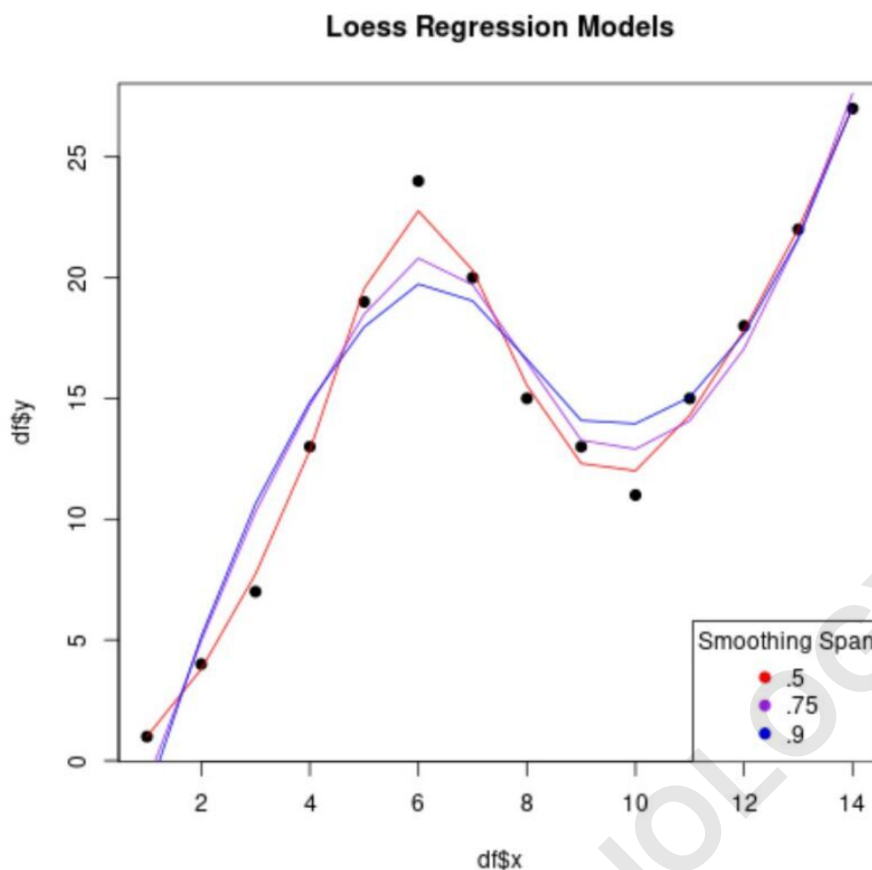
```
lines(smooth50, x=df$x, col='red')
```

```
lines(smooth75, x=df$x, col='purple')
```

```
lines(smooth90, x=df$x, col='blue')
```

```
legend('bottomright', legend=c('.5', '.75', '.9'),
```

```
col=c('red', 'purple', 'blue'), pch=19, title='Smoothing Span')
```



### Visualizing the Impact of Different Span Values

Upon reviewing the generated scatterplot with the overlaid LOESS curves, a critical pattern emerges concerning the smoothing **span**. The line corresponding to the smallest span (0.5, red line) is visibly less "smooth" and appears to oscillate closely around the individual data points. This behavior demonstrates high variance; the model is highly sensitive to the specific values of the local data subsets. While this tight fit minimizes the bias and captures rapid local changes, it increases the risk of overfitting, meaning the model may perform poorly when predicting new, unseen data.

Conversely, the line corresponding to the largest span (0.9, blue line) exhibits the highest degree of smoothness. This curve provides a much broader, more generalized view of the trend, minimizing variance but potentially introducing significant bias by smoothing over important local features in the dataset. This smooth fit is characteristic of underfitting, where the model is too rigid to capture the true complexity of the relationship.

The intermediate span (0.75, purple line) typically represents an attempt to strike a balance between these extremes. The objective is to find a curve that is smooth enough to generalize

beyond the immediate data points yet flexible enough to accurately represent the underlying non-linear pattern. Since visual inspection is subjective and unreliable for optimal model selection, a more rigorous, quantitative method is required to determine the best span value.

### Step 3: Use K-Fold Cross Validation to Find the Best Model

To move beyond subjective visual assessment and objectively identify the optimal **span** value, we must employ a robust model selection strategy. K-Fold Cross Validation is the gold standard for hyperparameter tuning, as it systematically assesses how well a model generalizes to independent data. This technique involves partitioning the dataset into K equally sized folds. The model is trained on K-1 folds and validated on the remaining fold. This process is repeated K times, ensuring every data point is used exactly once for validation.

In R, performing this advanced validation for LOESS models is typically achieved using the powerful `caret` package (Classification and Regression Training), which streamlines the process of defining training control parameters and iterating through candidate models. The method we employ here is 5-fold cross-validation (`cv`), meaning the data is divided into five parts. This rigorous testing procedure helps us select the span that consistently yields the lowest prediction error across various partitions of the data, thereby maximizing generalizability.

We define a grid of potential span values, ranging from 0.5 to 0.9, allowing the cross-validation framework to test each candidate. The goal is to find which span value minimizes the selected performance metric, typically the Root Mean Square Error (RMSE). This structured approach guarantees that the final model is selected based on its predictive performance rather than its fit to the training data alone.

### Executing Cross-Validation with the `caret` Package

To execute the cross-validation procedure, we first load the `caret` library and define the training control parameters. We specify a 5-fold cross-validation (`method = "cv", number = 5`) to ensure robust validation. Next, we define the search space for our hyperparameter--the `span`--using `expand.grid`, setting five equally spaced candidates between 0.5 and 0.9. We hold the polynomial degree constant at 1, which is common for initial LOESS modeling.

The `train` function within `caret` is then utilized to perform the cross-validation. We specify the model type as `"gamLoess"`, which allows `caret` to interface seamlessly with the LOESS fitting capabilities while managing the tuning process across the defined grid. This step automates the repetitive process of fitting the model with each span value, calculating the error on the validation fold, and averaging these errors across all five folds.

The resulting `model` object contains the comprehensive results of the cross-validation process,

including performance metrics for every tested span value. Printing this object provides a clear summary table, enabling a direct comparison of models based on their average predictive accuracy metrics, specifically focusing on the Root Mean Square Error (RMSE). This structured output is the definitive tool for optimal hyperparameter selection.

### library(caret)

```
#define k-fold cross validation method
ctrl <- trainControl(method = "cv", number = 5)
grid <- expand.grid(span = seq(0.5, 0.9, len = 5), degree = 1)

#perform cross-validation using smoothing spans ranging from 0.5 to 0.9
model <- train(y ~ x, data = df, method = "gamLoess", tuneGrid=grid, trControl = ctrl)
```

```
#print results of k-fold cross-validation
```

```
print(model)
```

```
14 samples
```

```
1 predictor
```

```
No pre-processing
```

```
Resampling: Cross-Validated (5 fold)
```

```
Summary of sample sizes: 12, 11, 11, 11, 11
```

```
Resampling results across tuning parameters:
```

```
span RMSE Rsquared MAE
```

```
0.5 10.148315 0.9570137 6.467066
```

```
0.6 7.854113 0.9350278 5.343473
```

```
0.7 6.113610 0.8150066 4.769545
```

```
0.8 17.814105 0.8202561 11.875943
```

```
0.9 26.705626 0.7384931 17.304833
```

```
Tuning parameter 'degree' was held constant at a value of 1
```

```
RMSE was used to select the optimal model using the smallest value.
```

```
The final values used for the model were span = 0.7 and degree = 1.
```

## Interpreting Results and Selecting the Optimal LOESS Model

The printed output from the cross-validation provides several key performance indicators for each tested span value. The primary metric used for selection in regression problems is the **Root Mean Square Error (RMSE)**, which quantifies the average magnitude of the errors in prediction. A lower

RMSE indicates a more accurate and generally superior predictive model. Other metrics, such as R-squared and Mean Absolute Error (MAE), also offer valuable context, but RMSE is the standard optimization target.

By examining the table, we observe a clear trend in the RMSE values. As the span increases from 0.5, the RMSE initially decreases, reaching its minimum at a span of **0.7** (RMSE = 6.113610). After this point, further increases in the span (0.8 and 0.9) cause the RMSE to sharply increase again, signifying a decline in model performance due to over-smoothing (underfitting).

Based on the quantitative evidence provided by the 5-fold cross-validation, we can confidently conclude that the value for **span** that produced the lowest RMSE is **0.7**. This span represents the ideal balance between minimizing bias and minimizing variance for this specific dataset. Thus, for our final LOESS regression model, we would choose to use a value of **0.7** for the **span** argument within the `loess()` function, ensuring the most accurate and generalizable fit.

## Summary of Best Practices for Local Regression

Successfully implementing LOESS regression involves more than just running the function; it requires careful consideration of the context and the critical parameter selection process. Always start with a visual inspection of the data to confirm that a non-linear method is appropriate. While LOESS is powerful, it is resource-intensive compared to global models and is best reserved for situations where its flexibility is necessary.

The determination of the optimal **span** is paramount. Never rely solely on visual fitting; always employ formal statistical techniques like K-Fold Cross Validation to objectively minimize prediction error metrics such as RMSE. Using robust packages like `caret` simplifies this iterative testing process, providing statistically sound evidence for the chosen model structure.

The following tutorials provide additional information about regression models in R and related data analysis techniques, offering pathways for further exploration of advanced statistical modeling methods: