

# How to Easily Generate Predictions from Linear Models in R Using `predict()` and `lm()`

Authored by  
**stats writer**

November 20, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Easily Generate Predictions from Linear Models in R Using `predict()` and `lm()`*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=98483>

## The Power of predict() with lm() in R

The ability to forecast future outcomes based on established statistical relationships is central to data analysis. In the [R programming language](#), this task is efficiently handled by coupling the **lm()** (linear model) function with the **predict()** function. The [linear model](#) serves as the mathematical foundation, quantifying the relationship between a response variable and one or more predictor variables. Once this model is fitted, the [predict\(\) function](#) is the essential tool for generating expected values.

The primary purpose of [predict\(\)](#) is to take a statistical model object--in this case, the output generated by **lm()**--and apply it to new, unseen data points, providing estimations of the outcome variable. This capability is vital for model validation, simulation studies, and deploying models into real-world applications where forecasting is required. It allows analysts to transition from descriptive statistics to predictive analytics seamlessly within the R environment.

Beyond simple point estimates, the **predict()** function offers crucial functionality for quantifying uncertainty. By utilizing optional arguments, users can generate both [confidence intervals](#) and prediction intervals for the resulting forecasts. Understanding how to correctly structure the input data and interpret the output of this function is fundamental for reliable statistical prediction in R.

The **lm()** function in R is specifically designed to fit [linear regression models](#).

Once we've fit a model using the observed data, we can then employ the [predict\(\) function](#) to calculate the estimated response value of a new observation. This process is commonly referred to as scoring the model.

### Syntax and Core Arguments of predict()

To effectively leverage the prediction capabilities within R, it is necessary to understand the standard syntax used by the function when applied to linear models. While the function has many optional parameters, a basic call requires the model object and, usually, a new data structure containing the input features for prediction.

This function uses the following general structure:

```
predict(object, newdata, type="response", interval="none", level=0.95)
```

The core parameters dictate how the prediction is executed and what auxiliary information is returned. Mastering these arguments is key to tailoring the prediction output to specific analytical needs.

The essential arguments are defined as follows:

**object:** This is the resulting object (typically of class 'lm') generated by fitting the linear model using the **lm()** function. This object contains all the necessary coefficients and model parameters required for the prediction calculation.

**newdata:** This highly crucial argument specifies a data frame containing the values of the predictor variables at which predictions are required. If **newdata** is omitted, the function defaults to predicting the values used to fit the original model.

**type:** This determines the type of prediction output. For **lm()** models, the default is "**response**", which provides the predicted mean value of the response variable. This is the output most frequently utilized in standard linear regression.

### Practical Demonstration: Setting Up the Data

To illustrate the application of predict() with **lm()**, we will use a dataset detailing basketball player statistics. Our goal is to predict the number of points a player scores based on minutes played and the number of fouls committed. This requires defining a specific dataset in R as a data frame.

Suppose we generate the following sample data frame in R, which captures information across ten different basketball players. This step ensures we have the necessary input variables (predictors) and output variable (response) before modeling begins.

**#create data frame**

```
df <- data.frame(minutes=c(5, 10, 13, 14, 20, 22, 26, 34, 38, 40),  
fouls=c(5, 5, 3, 4, 2, 1, 3, 2, 1, 1),  
points=c(6, 8, 8, 7, 14, 10, 22, 24, 28, 30))
```

**#view data frame**

```
df
```

```
minutes fouls points
```

```
1 5 5 6
```

```
2 10 5 8
```

```
3 1 3 8
```

```
4 14 4 7
```

```
5 20 2 14
```

```
6 22 1 10
```

```
7 26 3 22
```

```
8 34 2 24
```

```
9 38 1 28
```

```
10 40 1 30
```

This data frame, named **df**, contains our response variable, **points**, and our two predictor variables, **minutes** and **fouls**. The next step involves specifying the structure of the linear model we wish to estimate from this data.

## Fitting the Multiple Linear Regression Model

We aim to model the relationship where **points** is a linear function of **minutes** and **fouls**. This is expressed mathematically as a multiple linear regression equation:

$$\text{points} = \beta_0 + \beta_1(\text{minutes}) + \beta_2(\text{fouls}) + \varepsilon$$

Here,  $\beta_0$  represents the intercept, and  $\beta_1$  and  $\beta_2$  are the coefficients associated with minutes and fouls, respectively. We use the **lm()** function to estimate the optimal values for these coefficients that best fit our observed data, thereby minimizing the sum of squared residuals.

We execute the model fitting using **lm()** and store the result in an object named **fit**:

```
#fit multiple linear regression model
fit <- lm(points ~ minutes + fouls, data=df)
```

```
#view summary of model
summary(fit)
```

Call:

```
lm(formula = points ~ minutes + fouls, data = df)
```

Residuals:

```
Min 1Q Median 3Q Max
-3.5241 -1.4782 0.5918 1.6073 2.0889
```

Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
(Intercept) -11.8949 4.5375 -2.621 0.0343 *
minutes 0.9774 0.1086 9.000 4.26e-05 ***
fouls 2.1838 0.8398 2.600 0.0354 *
```

---

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

Residual standard error: 2.148 on 7 degrees of freedom

Multiple R-squared: 0.959, Adjusted R-squared: 0.9473

F-statistic: 81.93 on 2 and 7 DF, p-value: 1.392e-05

The summary output confirms the strength of the relationship, highlighted by the high Multiple R-squared value of 0.959. Using the coefficients estimated by **lm()**, we can formally write the derived regression equation that forms the basis of our predictions:

$$\text{points} = -11.8949 + 0.9774(\text{minutes}) + 2.1838(\text{fouls})$$

This equation is now parameterized and ready to be used by the **predict()** function to estimate outcomes for players not included in the original training data.

## Generating Predictions for New Observations

The true utility of a regression model lies in its predictive capability. We now apply the fitted model object (**fit**) to a hypothetical new player using the **predict()** function. Suppose we want to determine the expected points for a player who logs 15 minutes and commits 3 fouls.

To perform this prediction, we must construct a **newdata** data frame containing this specific observation. Crucially, the column names in this new data frame must exactly match the predictor names ('minutes' and 'fouls') used in the original **lm()** call.

We define the new observation and then pass it to the **predict()** function along with the fitted model object:

```
#define new observation  
newdata = data.frame(minutes=15, fouls=3)
```

```
#use model to predict points value  
predict(fit, newdata)
```

```
1  
9.317731
```

The output indicates that the model predicts this player, with 15 minutes played and 3 total fouls, will score **9.317731** points. This result represents the point estimate of the mean predicted response based on the linear relationship established by the training data.

## Predicting Multiple Observations Simultaneously

The **predict()** function is highly versatile and optimized to handle batch predictions efficiently. Instead of predicting one observation at a time, we can provide a **newdata** data frame containing multiple rows, where each row represents a distinct new observation for which we require a prediction.

For instance, we can predict the points scored for three different hypothetical players: Player A (15 minutes, 3 fouls), Player B (20 minutes, 2 fouls), and Player C (25 minutes, 1 foul). This scenario showcases how the **predict()** function handles vectorized inputs seamlessly.

The following code demonstrates how to define the batch of new observations and execute the prediction using the fitted **lm()** object:

```
#define new data frame of three players  
newdata = data.frame(minutes=c(15, 20, 25),  
fouls=c(3, 2, 1))
```

```
#view data frame  
newdata
```

```
minutes fouls  
1 15 3  
2 20 2  
3 25 1
```

```
#use model to predict points for all three players  
predict(fit, newdata)
```

```
1 2 3  
9.317731 12.021032 14.724334
```

The result is a vector of predicted values, corresponding positionally to the rows in the **newdata** data frame. Here is the interpretation of the output for each hypothetical player, rounding to two decimal places for practical use:

The predicted points for the player with 15 minutes and 3 fouls is **9.32**.

The predicted points for the player with 20 minutes and 2 fouls is **12.02**.

The predicted points for the player with 25 minutes and 1 foul is **14.72**.

## Advanced Usage: Incorporating Prediction and Confidence Intervals

While point predictions provide an estimated mean value, statistical rigor often demands quantifying the uncertainty around these estimates. The **predict()** function facilitates this by allowing the computation of confidence intervals and prediction intervals using the optional **interval** argument.

A confidence interval quantifies the uncertainty in the estimate of the **mean response** for a given set of predictors. These intervals are suitable when estimating the average outcome for a large

group of individuals who share the same characteristics. We compute this by setting `interval="confidence"`, and the confidence level can be adjusted using the `level` argument (default is 0.95).

In contrast, a prediction interval quantifies the uncertainty around a **single future observation**. Because single observations include the model's residual error, prediction intervals are inherently wider than confidence intervals for the same input values. This is the preferred interval type when forecasting an individual outcome. To calculate this, we set `interval="prediction"`.

## Critical Note on Data Structure: The newdata Requirement

One of the most frequent sources of errors when using the predict() function is a structural mismatch between the **newdata data frame** and the structure of the data used to train the original **lm()** model. Strict adherence to variable naming, data types, and ordering is mandatory.

Specifically, the names of the columns in the new data frame must exactly match the names of the predictor variables used in the formula passed to **lm()**. In our earlier example, the model fitting relied on the predictor variables:

**minutes**

**fouls**

Thus, when we created the new data frame called **newdata** we made sure to also name the columns:

**minutes**

**fouls**

If the column names are misspelled, miscapitalized, or if a required variable is missing, R will fail to evaluate the model equation correctly. This frequently results in a cryptic error message:

### Error in eval(predvars, data, env)

This error message indicates that R cannot find the variables it needs within the provided environment (the **newdata** object) to apply the model coefficients. Always keep this structural requirement in mind when preparing data for prediction using the **predict()** function.

## Summary of Key Predict() Concepts

The **predict()** function is an indispensable component of statistical modeling in R, especially when working with objects created by **lm()**. It provides the mechanism to extrapolate the findings of the

training data onto new, unobserved inputs.

Successful implementation relies on three key principles: first, ensuring the model object is correctly fitted using **lm()**; second, providing the prediction input via the **newdata** argument, structured identically to the training data; and third, selecting the appropriate prediction type (e.g., simple response, confidence interval, or prediction interval). Mastering these elements ensures powerful and reliable forecasting capabilities within the R environment.

ARABPSYCHOLOGY.COM