

Understanding Stepwise Regression: A Practical Guide with R Examples

Authored by
Mohammed loot

November 9, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Understanding Stepwise Regression: A Practical Guide with R Examples*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=14225>

The methodology of [Stepwise regression](#) provides an automated approach for constructing an optimal statistical [regression model](#). This procedure systematically selects or eliminates potential [predictor variables](#) from a larger set based on statistical criteria, such as minimizing the Akaike Information Criterion (AIC). The process iterates, adding or removing predictors sequentially until a statistically sound and parsimonious model is achieved, ensuring that only variables contributing meaningfully to the prediction remain.

The primary objective when employing stepwise selection is to refine a [multiple linear regression](#) model, ensuring it incorporates only those variables that demonstrate a statistically significant relationship with the desired [response variable](#). This balance between complexity and predictive power is crucial, particularly when dealing with large datasets where many correlated predictors exist, making manual variable selection cumbersome or unreliable.

Defining the Three Modes of Stepwise Selection

Stepwise regression is not a single process but a family of algorithms. Each approach utilizes a different strategy for navigating the vast space of possible models that can be constructed from a given set of predictor variables. Understanding these differences is essential for interpreting the results and selecting the most appropriate method for a specific analytical goal.

This comprehensive guide is designed to walk you through the practical application of stepwise regression using the statistical programming language [R](#). We will demonstrate how to implement the three foundational types of stepwise procedures, each offering a distinct approach to variable selection:

Forward Stepwise Selection: Beginning with an empty model and sequentially adding the most significant predictors one at a time.

Backward Stepwise Selection: Starting with the full model (all predictors included) and sequentially removing the least significant predictors.

Both-Direction Stepwise Selection: A hybrid approach that allows predictors to be added and subsequently removed in the search for the optimal model fit, providing the flexibility of both forward and backward movements.

Before diving into the R implementation, we must define our dataset, the core variables, and the specific function we will use to automate the selection process.

Data Setup and the R `step()` Function

To illustrate these methods practically, we will utilize the renowned built-in `mtcars` dataset in R, which provides data on 32 automobiles. This dataset is ideal for demonstrating model selection, as it contains numerous potential predictor variables relative to the sample size. We begin by

examining the structure of the data:

```
#view first six rows of mtcars
```

```
head(mtcars)
```

```
mpg cyl disp hp drat wt  qsec vs am gear carb
Mazda RX4 21.0 6 160 110 3.90 2.620 16.46 0 1 4 4
Mazda RX4 Wag 21.0 6 160 110 3.90 2.875 17.02 0 1 4 4
Datsun 710 22.8 4 108 93 3.85 2.320 18.61 1 1 4 1
Hornet 4 Drive 21.4 6 258 110 3.08 3.215 19.44 1 0 3 1
Hornet Sportabout 18.7 8 360 175 3.15 3.440 17.02 0 0 3 2
Valiant 18.1 6 225 105 2.76 3.460 20.22 1 0 3 1
```

Our objective is to construct a robust [multiple linear regression](#) model where miles per gallon (*mpg*) serves as the primary **response variable**. The remaining ten variables, ranging from weight (*wt*) to horsepower (*hp*) and number of cylinders (*cyl*), will be considered as the full set of potential **predictor variables** in our model selection process.

All three forms of stepwise selection in R are handled efficiently by the built-in [step\(\)](#) function, which is part of the standard `stats` package. This function utilizes the [AIC](#) (Akaike Information Criterion) as the primary metric for judging model quality, seeking to minimize this value. The general syntax requires three key arguments:

```
step(initial_model, direction, scope)
```

The parameters define the starting point and limits of the search process:

initial_model: This is the starting point for the search. For forward selection, this will typically be the minimal model (intercept-only). For backward selection, this is usually the maximal model (including all predictors).

direction: This critical argument dictates the mode of search, specified as `"forward"`, `"backward"`, or `"both"`.

scope: This defines the range of models the search can explore. It usually specifies the maximal model (upper bound, containing all candidate variables) to ensure the search is constrained to the variables of interest.

The stepwise algorithm proceeds by comparing the AIC of potential models at each stage. A lower AIC indicates a better trade-off between the model's goodness of fit and the complexity introduced by additional **predictor variables**.

Example 1: Implementing Forward Stepwise Selection in R

Forward selection begins with the simplest possible model--the **intercept-only model**--which contains no predictor variables. The algorithm then iteratively tests adding each available predictor one by one. At each step, the variable that produces the largest statistically significant decrease in the [AIC](#) is added to the model. The process halts when no remaining variable can significantly improve the AIC score, ensuring the final model is as parsimonious as possible while maximizing explanatory power.

The following R code defines both the minimal intercept model and the maximal model containing all predictors, then executes the forward stepwise search using `direction='forward'`. Note that we define the scope using the `formula(all)` object, which lists all candidate variables available for inclusion.

```
#define intercept-only model (starting point)
```

```
intercept_only <- lm(mpg ~ 1, data=mtcars)
```

```
#define model with all predictors (defines the scope)
```

```
all <- lm(mpg ~ ., data=mtcars)
```

```
#perform forward stepwise regression (trace=0 suppresses verbose output)
```

```
forward <- step(intercept_only, direction='forward', scope=formula(all), trace=0)
```

```
#view results of forward stepwise regression
```

```
forward$anova
```

```
Step Df Deviance Resid. Df Resid. Dev AIC
1 NA NA 31 1126.0472 115.94345
2 + wt -1 847.72525 30 278.3219 73.21736
3 + cyl -1 87.14997 29 191.1720 63.19800
4 + hp -1 14.55145 28 176.6205 62.66456
```

```
#view final model coefficients
```

```
forward$coefficients
```

```
(Intercept) wt cyl hp
38.7517874 -3.1669731 -0.9416168 -0.0180381
```

The `trace=0` argument is used here to ensure that R only displays the final summary table (`$anova`) rather than showing the results of every variable permutation tested, which can quickly become overwhelming when working with many **predictor variables**. The interpretation of the

`$anova` output details the sequence of variable inclusion and the corresponding reduction in the AIC:

Step 1 (Baseline): The process starts with the intercept-only model, which has a high baseline AIC of **115.94345**.

Step 2 (+ wt): The variable *wt* (weight) provided the most significant improvement in fit, reducing the AIC substantially to **73.21736**.

Step 3 (+ cyl): With *wt* already in the model, *cyl* (number of cylinders) was identified as the next best predictor, further lowering the AIC to **63.19800**.

Step 4 (+ hp): The inclusion of *hp* (horsepower) resulted in a marginal but significant improvement, yielding a final AIC of **62.66456**.

Stopping Criterion: After adding *hp*, the algorithm determined that no other available predictor could significantly decrease the AIC. Consequently, the procedure stopped, identifying the three variables (*wt*, *cyl*, *hp*) as the optimal set.

Based on the coefficients of the final model, the resulting regression equation derived from the forward selection process is:

$$\text{mpg} = 38.75 - 3.17 * \text{wt} - 0.94 * \text{cyl} - 0.02 * \text{hp}$$

Example 2: Executing Backward Stepwise Selection

In contrast to the forward method, **Backward stepwise selection** begins with the **maximal model**, which includes all potential [predictor variables](#) (*p* predictors). The algorithm then systematically removes the least statistically significant predictor one at a time. A variable is dropped if its removal results in the smallest increase (or ideally, a decrease) in the [AIC](#), suggesting that the variable was unnecessary noise or redundancy. This iterative elimination continues until the removal of any remaining variable would cause a significant deterioration in the model fit (a substantial increase in AIC).

For backward selection, the `step()` function is initialized with the full model (`all`) and the direction is set to `'backward'`. This method is often preferred by statisticians because it forces the algorithm to consider all possible interactions and confounding effects present in the full model before making elimination decisions.

#define intercept-only model (required for scope definition)

```
intercept_only <- lm(mpg ~ 1, data=mtcars)
```

#define model with all predictors (starting point)

```
all <- lm(mpg ~ ., data=mtcars)
```

```
#perform backward stepwise regression
```

```
backward <- step(all, direction='backward', scope=formula(all), trace=0)
```

```
#view results of backward stepwise regression
```

```
backward$anova
```

```
Step Df Deviance Resid. Df Resid. Dev AIC
1 NA NA 21 147.4944 70.89774
2 - cyl 1 0.07987121 22 147.5743 68.91507
3 - vs 1 0.26852280 23 147.8428 66.97324
4 - carb 1 0.68546077 24 148.5283 65.12126
5 - gear 1 1.56497053 25 150.0933 63.45667
6 - drat 1 3.34455117 26 153.4378 62.16190
7 - disp 1 6.62865369 27 160.0665 61.51530
8 - hp 1 9.21946935 28 169.2859 61.30730
```

```
#view final model
```

```
backward$coefficients
```

```
(Intercept) wt qsec am
9.617781 -3.916504 1.225886 2.935837
```

The results show that the process began with the full model (AIC 70.89774) and sequentially removed seven variables (*cyl*, *vs*, *carb*, *gear*, *drat*, *disp*, *hp*) because their exclusion either marginally increased or decreased the overall AIC, indicating they were not essential to the model's predictive accuracy. The final model selected three **predictor variables**: *wt*, *qsec*, and *am*, resulting in the lowest overall AIC of 61.30730. This demonstrates a key difference from the forward selection method, which yielded a set of predictors (*wt*, *cyl*, *hp*) and a slightly higher AIC (62.66456).

The final refined model based on the backward elimination procedure is formulated as:

$$\text{mpg} = 9.62 - 3.92 * \text{wt} + 1.23 * \text{qsec} + 2.94 * \text{am}$$

Example 3: Utilizing Both-Direction (Hybrid) Stepwise Selection

The **Both-Direction** or **Hybrid stepwise selection** combines the search mechanisms of both forward addition and backward elimination into a single, comprehensive strategy. Starting typically from the intercept-only model, it adds the best predictor at each stage (like forward selection), but then critically re-examines all variables currently in the model. If a variable that was previously added no longer contributes significantly (based on the [AIC](#)), it is removed before the next addition step. This is particularly useful for mitigating the effects of [multicollinearity](#), where adding one

variable might render a previously included variable redundant.

To implement this hybrid method, we set the starting model to the intercept-only model (`intercept_only`) and specify the direction argument as `'both'`. The algorithm will then determine whether to add a new variable or remove an existing one in its attempt to minimize the AIC.

```
#define intercept-only model (starting point)
```

```
intercept_only <- lm(mpg ~ 1, data=mtcars)
```

```
#define model with all predictors (defines the scope)
```

```
all <- lm(mpg ~ ., data=mtcars)
```

```
#perform both-direction stepwise regression
```

```
both <- step(intercept_only, direction='both', scope=formula(all), trace=0)
```

```
#view results of both-direction stepwise regression
```

```
both$anova
```

```
Step Df Deviance Resid. Df Resid. Dev AIC
```

```
1 NA NA 31 1126.0472 115.94345
```

```
2 + wt -1 847.72525 30 278.3219 73.21736
```

```
3 + cyl -1 87.14997 29 191.1720 63.19800
```

```
4 + hp -1 14.55145 28 176.6205 62.66456
```

```
#view final model
```

```
both$coefficients
```

```
(Intercept) wt cyl hp
```

```
38.7517874 -3.1669731 -0.9416168 -0.0180381
```

The interpretation of the results aligns closely with the forward selection example because, in this specific application of the **mtcars** dataset, no variables that were added early on became redundant or warranted removal later in the process. The sequence of variable addition (*wt*, *cyl*, *hp*) mirrors the forward selection results perfectly, resulting in an identical final model and AIC score of **62.66456**.

The final model obtained through the both-direction search is:

```
mpg = 38.75 - 3.17 * wt - 0.94 * cyl - 0.02 * hp
```

Comparing Stepwise Selection Methods and Key Takeaways

A critical observation from these examples is that different selection methodologies can lead to divergent final models, even when applied to the exact same dataset and objective. In our analysis, the forward selection and the both-direction selection methods converged on the model containing *wt*, *cyl*, and *hp*. Conversely, the backward elimination approach selected a different set: *wt*, *qsec*, and *am*, which yielded a marginally superior AIC (61.30730) compared to the forward/both model (62.66456).

This difference highlights a potential weakness of stepwise methods: they may settle on a local optimum rather than the global optimal model. Forward selection, for instance, might exclude a variable early on simply because it doesn't meet the inclusion threshold in isolation, even if it would become highly significant when paired with another variable added later. Since backward elimination considers the full model from the start, it often provides a more robust starting point, though it is computationally more intensive for datasets with an exceptionally large number of **predictor variables**.

Regardless of the method chosen, it is essential to remember that stepwise regression is a tool for exploration and refinement, not definitive proof of causality. The resulting models should always be subjected to rigorous external validation, and evaluated for compliance with the underlying assumptions of [multiple linear regression](#) (such as linearity, homoscedasticity, and independence of errors) before being adopted for predictive purposes.

Additional Resources

[How to Test the Significance of a Regression Slope](#)

[How to Read and Interpret a Regression Table](#)

[A Guide to Multicollinearity in Regression](#)