

Learning K-Fold Cross-Validation with R: A Step-by-Step Guide

Authored by
Mohammed loot

November 6, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning K-Fold Cross-Validation with R: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=11865>

To properly assess the performance of a predictive [model](#) on new, unseen data, it is essential to employ robust evaluation techniques. When developing machine learning or statistical models, the primary goal is often generalization--the ability of the model to make accurate predictions on data it was not trained on. A simple train/test split can sometimes lead to overly optimistic or pessimistic results, especially with smaller datasets.

To overcome the limitations of the standard holdout method and obtain a more reliable estimate of the model's prediction error, data scientists commonly utilize [cross-validation](#). The most widely adopted form of this technique is known as **K-Fold Cross-Validation**, a powerful method designed to measure how well the model's predictions align with the observed data across the entire dataset.

Understanding K-Fold Cross-Validation Methodology

K-Fold Cross-Validation is an iterative process that systematically rotates validation and training subsets, ensuring every observation in the dataset has the opportunity to be included in the [test set](#) exactly once. This procedure yields a less biased and more stable estimate of model performance compared to a single data split.

The methodology relies on dividing the input data into a specified number of sections, denoted by k . While the choice of k is flexible, common values are 5 or 10. The core steps of this approach are as follows:

The entire dataset is randomly partitioned into k equal-sized subgroups, or "folds."

In the first iteration, one fold is designated as the validation or holdout set (for testing), and the model is fitted (trained) using the remaining $k-1$ folds. After training, the performance metric (such as the Mean Squared Error, or MSE) is calculated specifically on the observations in the held-out fold.

This training and validation process is repeated k times. In each subsequent iteration, a different one of the k folds is used as the validation set.

Finally, the overall model performance metric, often the overall test MSE or [Root Mean Squared Error \(RMSE\)](#), is determined by calculating the average of the k individual test performance metrics collected from each iteration.

This averaging process is critical because it mitigates the randomness associated with a single training/testing split, providing a final performance score that reflects the model's robustness across different subsets of the data.

Implementing K-Fold Cross-Validation in R

The easiest and most efficient way to execute K-Fold Cross-Validation in the [R programming](#)

[language](#) is by leveraging the powerful **caret** package (Classification And REgression Training). This package offers a unified interface for various machine learning tasks, including model training, tuning, and resampling methods.

Specifically, the `trainControl()` function within the [caret package](#) allows us to define the resampling technique we wish to use, such as cross-validation, and specify the number of folds (k). The subsequent `train()` function then handles the model fitting and iterative cross-validation process automatically based on these control settings.

This tutorial provides a practical, step-by-step demonstration of how to implement 5-fold cross-validation for a standard regression model using these functions in R.

Example: Setting Up the Dataset and Model

To illustrate the process, let us define a small example dataset in R. This dataset contains ten observations, with a response variable y and two predictor variables, x_1 and x_2 .

```
#create data frame
df <- data.frame(y=c(6, 8, 12, 14, 14, 15, 17, 22, 24, 23),
x1=c(2, 5, 4, 3, 4, 6, 7, 5, 8, 9),
x2=c(14, 12, 12, 13, 7, 8, 7, 4, 6, 5))

#view data frame
df

y x1 x2
6 2 14
8 5 12
12 4 12
14 3 13
14 4 7
15 6 8
17 7 7
22 5 4
24 8 6
23 9 5
```

Our objective is to fit a [multiple linear regression model](#) predicting y based on x_1 and x_2 , and then evaluate its stability and predictive accuracy using 5-fold cross-validation. This ensures that the evaluation metrics we obtain are representative of the model's performance on future data.

The following R code demonstrates the necessary steps: first loading the `caret` library, then defining the cross-validation parameters using `trainControl()`, and finally training the regression model using the `train()` function while applying the defined CV controls.

library(caret)

```
#specify the cross-validation method
ctrl <- trainControl(method = "cv", number = 5)

#fit a regression model and use k-fold CV to evaluate performance
model <- train(y ~ x1 + x2, data = df, method = "lm", trControl = ctrl)

#view summary of k-fold CV
print(model)
```

Linear Regression

10 samples
2 predictor

No pre-processing
Resampling: Cross-Validated (5 fold)
Summary of sample sizes: 8, 8, 8, 8, 8
Resampling results:

RMSE Rsquared MAE
3.018979 1 2.882348

Tuning parameter 'intercept' was held constant at a value of TRUE

Interpreting the Cross-Validation Output Metrics

The summary output generated by the `print(model)` command provides essential information regarding the cross-validation setup and the resulting average performance metrics. Understanding these metrics is crucial for determining the quality of the regression model.

The output confirms that our setup used 5-fold cross-validation (Resampling: Cross-Validated (5 fold)), and since the total sample size is 10, each training set consisted of 8 observations (10 total - 2 in the test fold). The key figures are summarized in the Resampling results table:

RMSE (Root Mean Squared Error): This is perhaps the most widely used measure of prediction accuracy for regression models. It quantifies the average magnitude of the errors, where the errors

are the differences between the model's predictions and the actual observations. Because errors are squared before being averaged, larger errors are penalized more heavily. A **lower RMSE** value indicates that the model's predictions closely align with the true data points, signifying better performance.

R-squared (Coefficient of Determination): This metric measures the proportion of the variance in the dependent variable that is predictable from the independent variables. R-squared ranges from 0 to 1 (or 0% to 100%). In the context of cross-validation, the reported R-squared is the average performance across all folds. A **higher R-squared** suggests that the model explains a greater proportion of the variability in the response, indicating a stronger relationship between the predictors and the outcome.

MAE (Mean Absolute Error): This metric calculates the average absolute difference between the predicted values and the actual observed values. Unlike RMSE, MAE is less sensitive to outliers because it does not square the errors. Like RMSE, a **lower MAE** indicates a higher degree of predictive accuracy.

These three metrics collectively provide a comprehensive view of how well the model generalizes to previously unseen data, based on the averaged results from the five validation folds. In applied statistics, practitioners often fit several competing models and use these metrics to objectively decide which model produces the lowest test error rates and is therefore the superior choice.

Examining the Final Model Equation and Fold Variability

While the summary output provides the averaged performance metrics, we can also extract the coefficients of the final model fitted on the entire dataset, after the cross-validation process is complete. This final model is the one typically deployed for future predictions.

```
#view final model  
model$finalModel
```

Call:

```
lm(formula = .outcome ~ ., data = dat)
```

Coefficients:

```
(Intercept) x1 x2
```

```
21.2672 0.7803 -1.1253
```

Based on the coefficients provided in the output, the final estimated linear regression equation is derived as:

$$y = 21.2672 + 0.7803 * (x1) - 1.1253 * (x2)$$

Furthermore, one of the significant benefits of using K-Fold Cross-Validation is the ability to inspect the performance metrics for each individual fold. This allows us to assess the stability and variance of the model's performance across different subsets of the data. High variability between fold metrics might suggest that the model is sensitive to the specific training data selected.

We can access this detailed breakdown using the `model$resample` component:

#view predictions for each fold

model\$resample

```
RMSE Rsquared MAE Resample
1 4.808773 1 3.544494 Fold1
2 3.464675 1 3.366812 Fold2
3 6.281255 1 6.280702 Fold3
4 3.759222 1 3.573883 Fold4
5 1.741127 1 1.679767 Fold5
```

Observing the results above, we can see the performance metrics (RMSE, R-squared, MAE) differ for each of the five folds. For instance, Fold 5 yielded a very low RMSE (1.74), suggesting excellent performance on that specific holdout set, whereas Fold 3 had a significantly higher RMSE (6.28). The overall RMSE reported earlier (3.018979) is the average of these five results, providing a balanced estimate of generalization error.

Considerations for Choosing the Value of K

While we utilized $k=5$ folds in this demonstration, the choice of k is a practical decision involving a trade-off between computational cost and bias/variance of the resulting error estimate.

If k is small (e.g., $k=2$), the model is trained on less data (only half the dataset), leading to a higher bias in the model fit, although the computation is faster. Conversely, if k is large (e.g., k equals the number of observations, known as Leave-One-Out Cross-Validation or LOOCV), the model is trained on almost all the data, leading to low bias but high variance in the error estimate and significantly increased computational time.

In practical machine learning and statistical modeling, the values **$k=5$ or $k=10$** are standard conventions. These numbers typically strike the optimal balance, producing reliable test error rates that are not overly biased while maintaining manageable computational expense, making K-Fold Cross-Validation a cornerstone technique in model evaluation.