

Use createDataPartition() Function in R

Authored by
Mohammed loot

March 19, 2026

RECOMMENDED CITATION

Mohammed loot (2026). *Use createDataPartition() Function in R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=3299>

In the realm of [machine learning](#), the meticulous preparation of data stands as a critical prerequisite that fundamentally dictates the performance, stability, and reliability of any subsequent [predictive model](#). A cornerstone of this preparation methodology involves the systematic division of the complete dataset into distinct, non-overlapping subsets intended for training and rigorous testing. This essential procedure is paramount for accurately evaluating how effectively a model can generalize its learned patterns to entirely unseen data, thereby serving as the primary defense against common and detrimental issues such as [overfitting](#). The [R programming language](#) furnishes data scientists with powerful and adaptable tools for this requirement, among which the **`createDataPartition()`** function, housed within the universally adopted [caret](#) package, is recognized as a highly flexible and robust solution.

The Indispensable Role of Data Partitioning in Machine Learning

Before initiating a deep dive into the technical specificities of **`createDataPartition()`**, it is vital to establish a clear understanding of why dataset partitioning is not merely an optional step, but an absolutely indispensable requirement in every professional [machine learning](#) workflow. When we engage in the construction of a [predictive model](#), our core objective extends far beyond achieving high accuracy on the specific data utilized for its initial creation. Fundamentally, we strive for a model that exhibits the capacity to accurately forecast outcomes when presented with novel, previously unobserved data points. If the model is both trained and subsequently validated using the identical dataset, there is a significant risk that it will merely memorize the training examples and noise, rather than truly learning the underlying, generalizable patterns. This critical failure mode is universally termed **overfitting**, and it invariably results in models that exhibit severely degraded performance when deployed in practical, real-world applications.

To effectively counteract the pervasive threat of overfitting and to provide an objective, unbiased assessment of a model's true generalization capability, the full dataset must be systematically separated into at least two distinct, primary components: the **training set** and the **testing set**. The **training set** is exclusively designated for teaching the model, enabling it to iteratively discern complex relationships, extract salient features, and identify patterns inherent within the data. Once this learning phase is complete, the **testing set**--which comprises data instances the model has categorically never encountered--is utilized to rigorously evaluate the model's aptitude for generalization. This stringent separation is the only mechanism that guarantees the evaluation metrics genuinely reflect the model's predictive capabilities on new, unobserved data, thereby offering the most reliable measure of its actual predictive power in deployment scenarios.

While rudimentary random sampling techniques can be employed to arbitrarily split datasets, this approach often carries the significant liability of inadvertently generating imbalanced distributions of the crucial target variable across the resulting training and testing subsets. This specific issue is exacerbated in two scenarios: when dealing with smaller datasets, or more commonly, when

working with datasets characterized by [imbalanced classes](#), where certain categories are fundamentally underrepresented. In these challenging contexts, sophisticated partitioning strategies, such as the capabilities provided by functions like `createDataPartition()`, transition from being useful to becoming absolutely invaluable. These advanced methods ensure that the essential statistical properties of the original complete dataset are meticulously preserved across both the training and testing subsets, a feature which is crucial for yielding robust, representative samples and ultimately, more reliable predictive outcomes.

Introducing `createDataPartition()` from the Comprehensive `caret` Package

The [caret](#) package, an acronym aptly derived from **C**lassification **A**nd **R**egression **T**raining, stands as an exceptionally comprehensive and highly versatile collection of functions integrated within the [R](#) environment. This package has been meticulously engineered to significantly streamline and simplify the entirety of the [machine learning](#) workflow. Its robust utilities span the full spectrum of development, including initial data preprocessing, crucial feature selection, rigorous model training, and thorough performance evaluation. Prominently featured among its extensive suite of functions, `createDataPartition()` is specifically designed to generate reliable and statistically representative data splits, which is recognized as a fundamental cornerstone of effective and responsible [predictive modeling](#).

The primary functional objective of `createDataPartition()` is to enable the creation of [stratified random samples](#). This highly sophisticated sampling technique rigorously guarantees that, upon partitioning the data, the proportion of various classes (a necessity for classification problems) or the overall distribution characteristics of the outcome variable (essential for regression tasks) within both the **training set** and the **testing set** closely and accurately mirrors those found in the original, complete dataset. This meticulous and careful preservation of the data's inherent statistical characteristics is critically important in contexts where certain classes are inherently rare, or where the outcome variable's natural distribution is significantly skewed. Without this process of stratification, an arbitrary random split could easily lead to one subset being statistically unrepresentative, thereby compromising the fundamental integrity of both the model training process and the subsequent evaluation results.

By intelligently and systematically employing [stratified sampling](#), `createDataPartition()` fulfills a vital role in facilitating the development of models that are not only more robust but also significantly more reliable. It ensures that the model learns from a statistically diverse and truly representative sample of the data, and crucially, is subsequently evaluated on a test set that accurately reflects the real-world data distribution it is ultimately expected to encounter. This disciplined and careful approach to the partitioning of data is instrumental in yielding performance assessments that are substantially more dependable, directly contributing to superior, trustworthy predictive capabilities in virtually all practical applications.

Mastering the Syntax and Essential Parameters

The `createDataPartition()` function is defined by an exceptionally clear and highly intuitive syntax, making it readily accessible for addressing a wide spectrum of data partitioning requirements encountered in data science. Its fundamental operational structure offers flexible yet precise control over the splitting procedure, empowering users to customize the process meticulously according to their distinct analytical needs and modeling objectives.

`createDataPartition(y, times = 1, p = 0.5, list = TRUE, ...)`

To fully leverage this powerful function, we must now thoroughly examine the role and specific contribution of each core parameter in the robust generation of data partitions:

y: This critically important parameter mandates a **vector of outcomes**. For [regression problems](#), the `y` input must be a numeric vector that accurately represents the continuous target variable. Conversely, for [classification problems](#), it should be a factor vector explicitly denoting the categorical classes. The values contained within `y` are the central input for the stratified sampling mechanism. The function utilizes this information to intelligently create balanced splits, thereby ensuring that the distribution profile of the outcome variable is meticulously and proportionally maintained across both the training and testing sets. In classification scenarios, this specifically entails preserving the approximate percentage of each class; in regression, the goal is to balance the spread and range of the numeric outcome values.

times: This integer parameter rigorously dictates the **number of distinct partitions to be generated**. A value of `1`, which serves as the function's default setting, signifies the creation of a single split, resulting in precisely one pair of training and testing sets. If a user specifies a value greater than `1` for `times`, the function proceeds to produce multiple, entirely independent training and testing splits. This specialized capability is particularly advantageous for implementing sophisticated [repeated cross-validation](#) procedures, which significantly enhance the statistical reliability of model evaluation by systematically averaging performance metrics across several distinct splits.

p: This numeric parameter accepts a fractional value ranging strictly from `0` to `1` and determines the precise **percentage of the total data to be allocated specifically to the training set**. For example, setting `p` to `0.8` (or simply `.8`) dictates that 80% of the complete dataset will be assigned to the training set, with the resultant 20% being reserved for the testing set. Commonly utilized values in professional [machine learning](#) practice often include `0.7`, `0.75`, or `0.8`, with the selection being carefully chosen based on critical factors such as the overall dataset size and the specific complexity requirements of the predictive model being developed.

list: This logical parameter controls the fundamental structure of the output, specifically determining **whether the results are stored within a list data structure**. When this parameter is

set to `TRUE` (which is the default setting), the function returns a list where each element is a vector containing the row indices designated for inclusion in the training set. Conversely, if `list` is explicitly set to `FALSE`, the function directly returns a single consolidated vector of row indices, a format that is often significantly more convenient when the user's requirement is to create only a singular data partition. It is critical to note a functional exception: for the generation of multiple partitions (i.e., when `times > 1`), the function will invariably return a list of vectors, even if `list = FALSE` was specified, with each element of the list containing a vector of indices corresponding to one training split.

...: The ellipsis symbol provides a mechanism for the passing of **additional arguments** to underlying functions that `createDataPartition()` might utilize internally during its execution. While its inclusion ensures maximum flexibility and extensibility, these additional arguments are generally not required or directly specified by users in the majority of standard, routine data partitioning scenarios.

Practical Implementation: Data Partitioning in R

To thoroughly demonstrate the practical utility and application of the `createDataPartition()` function, we will walk through a common analytical scenario where the objective is to construct a [linear regression model](#). Our process will commence by generating a realistic synthetic dataset, and subsequently, we will utilize the function to meticulously and strategically split this data into the required training and testing subsets. This comprehensive example is designed to clearly illustrate all necessary steps involved in preparing your data for highly effective model development and robust evaluation.

Preparing the Sample Data for Analysis

Let us construct a hypothetical [data frame](#) in [R](#) that is populated with 1,000 observations. This dataset is engineered to contain information tracking the **hours** students dedicated to studying and their corresponding final examination **score**. Our analytical aim is to build a [predictive model](#) capable of accurately estimating a student's examination score based purely on the observed number of hours they spent studying.

Prior to the critical step of data generation, it is recognized as an absolute best practice to establish and set a [random seed](#). This crucial action ensures the essential reproducibility of our entire example, guaranteeing that the random values generated will remain identical every single time the code block is executed, thereby facilitating consistent results, accurate model comparisons, and dramatically simpler debugging processes.

```
# Ensure the reproducibility of this example  
set.seed(0)
```

```
# Create the data frame containing hours and score variables
df <- data.frame(hours=runif(1000, min=0, max=10),
score=runif(1000, min=40, max=100))
```

```
# Display the first few rows of the data frame for verification
```

```
head(df)
```

```
hours score
1 8.966972 55.93220
2 2.655087 71.84853
3 3.721239 81.09165
4 5.728534 62.99700
5 9.082078 97.29928
6 2.016819 47.10139
```

This segment of code successfully generates a [data frame](#) designated as `df`. It is populated with two columns: `hours`, which contains random uniform values ranging between 0 and 10, and `score`, which is similarly populated with random uniform values ranging between 40 and 100. The application of the `head()` function is then employed to provide a rapid visual inspection of the initial rows of our newly constructed dataset, which serves to confirm its intended structure and the initial content generated.

Executing the Data Partition

Our primary objective remains the training of a [linear model](#) utilizing the `hours` variable to predict the continuous `score` outcome. To ensure the model evaluation is robust and fundamentally unbiased, we have strategically chosen to allocate a substantial 80% of our complete dataset specifically for the purpose of model training, reserving the remaining 20% for the critical step of performance testing. This specific 80/20 division ratio is a widely accepted and highly utilized standard ratio across diverse [machine learning](#) practices.

We will now proceed to utilize the powerful `createDataPartition()` function, which is readily available within the [caret](#) package, to execute this required stratified split. Given the nature of our outcome variable (`score`) being numeric and continuous, `createDataPartition()` will intelligently work to ensure that the distribution characteristics of the scores in both the training and testing sets closely and accurately mirror those of the original dataset. For enhanced convenience and simplicity in this specific demonstration scenario, we explicitly set the `list` parameter to `FALSE`, a command which instructs the function to return a direct vector composed solely of row indices, rather than the default list structure.

`library(caret)`

```
# Partition data frame into training and testing sets based on 'score'
train_indices <- createDataPartition(df$score, times=1, p=.8, list=FALSE)

# Create the training set using the generated row indices
df_train <- df

# Create the testing set by selecting rows that were NOT included in the training indices
df_test <- df

# View the definitive number of rows contained in each newly created set
nrow(df_train)

800

nrow(df_test)

200
```

The resulting `train_indices` object, which constitutes the primary output of the **`createDataPartition()`** function, is a vector containing the specific row numbers designated for mandatory inclusion in our **training set**. By systematically employing these generated indices, we successfully construct the `df_train` [data frame](#). Subsequently, the `df_test` [data frame](#) is constructed by strategically selecting all rows from the original `df` that were definitively *not* included in the training set allocation. This highly efficient inverse selection is achieved by placing a negative sign immediately before the `train_indices` vector, which effectively instructs R to select all rows *except* those explicitly specified in the vector, ensuring a clean and non-overlapping split.

Verifying and Validating the Partitioned Datasets

Following the execution of the data partitioning code, it is critically important to verify and confirm that the split procedure has been executed with the precision and ratio that was intended. We can effortlessly confirm the exact dimensions of our newly generated datasets by applying the fundamental [`nrow\(\)`](#) function, which returns the definitive count of rows contained within any given [data frame](#).

As was correctly anticipated and subsequently confirmed by the output displayed, our resulting **training set** (`df_train`) now meticulously contains exactly 800 rows, which precisely accounts for 80% of the original 1,000-row dataset. Concurrently, the complementary **testing set** (`df_test`) is composed of 200 rows, accurately representing the remaining 20% allocation. This verification step

unequivocally confirms that our desired 80/20 data split was executed successfully and with the required accuracy, providing us with distinct and properly sized datasets for the subsequent stages of model development and unbiased performance evaluation.

Furthermore, in order to visually ascertain the critical distinctiveness and non-overlapping nature of the two resultant datasets, we should inspect the initial rows of both the training and testing subsets by using the [head\(\)](#) function. This straightforward visual check serves as an important additional confirmation that the datasets are indeed entirely separate entities and contain different data points, each now appropriately prepared to fulfill its respective, designated role in the subsequent pipeline stages of model training and robust evaluation.

View the first few rows of the training set

[head\(df_train\)](#)

```
hours score
1 8.966972 55.93220
2 2.655087 71.84853
3 3.721239 81.09165
4 5.728534 62.99700
5 9.082078 97.29928
7 8.983897 42.34600
```

View the first few rows of the testing set

[head\(df_test\)](#)

```
hours score
6 2.016819 47.10139
12 2.059746 96.67170
18 7.176185 92.61150
23 2.121425 89.17611
24 6.516738 50.47970
25 1.255551 90.58483
```

Carefully observe that the row numbers explicitly presented in the training set (for instance, 1, 2, 3, 4, 5, 7) are definitively distinct and non-overlapping from those displayed in the testing set (such as 6, 12, 18, 23, 24, 25). This clear and demonstrable differentiation unequivocally confirms that the original data has been successfully and correctly partitioned into two entirely separate, mutually exclusive datasets. Consequently, each of these subsets is now fully and appropriately prepared for its designated, critical role in the model development and evaluation pipeline, thereby ensuring a scientifically rigorous and fair assessment of our [predictive model](#).

Why Stratified Sampling Ensures Robust Model Quality

One of the most consequential and compelling advantages that fundamentally separates `createDataPartition()` from any simple, rudimentary random split is its inherent, built-in capability to execute [stratified sampling](#). This core feature provides profound benefits, especially when analysts are confronted with datasets where the distribution of the target variable (which is specified via the mandatory `y` parameter) is uneven, or when certain critical categories are naturally rare or severely underrepresented within the overall population. By systematically maintaining the proportional representation of these crucial characteristics, the function guarantees that both the training and testing subsets are robustly and statistically representative of the original data's composition.

In the specific context of [classification tasks](#), for example, if an analyst is working with a highly imbalanced dataset where a minority class constitutes a very small fraction of the total observations, a purely random partitioning methodology might inadvertently lead to the creation of a [training set](#) or a [testing set](#) that completely lacks instances of this crucial minority class, or contains an insufficient number of them to be statistically meaningful. Such a pronounced imbalance can severely compromise the model's ability to effectively learn from and accurately predict that particular class. By leveraging the information contained in the `y` parameter, `createDataPartition()` meticulously ensures that the proportions of every single class are largely preserved across both the training and testing sets. This leads directly to more representative samples and, consequently, results in substantially more reliable and unbiased performance evaluations of the model.

Similarly, when dealing with [regression problems](#) which involve continuous outcomes, `createDataPartition()` actively strives to maintain the fidelity of the overall distribution of the continuous outcome variable across all generated splits. This highly sophisticated approach ensures that the range, density, and general spread characteristics of the `y` values found in the training data will closely and accurately mirror those present within the test data. This consistency provides a significantly more stable and comparable foundation for both the initial model training and the subsequent rigorous evaluation. This meticulous attention to the detail of sampling is recognized as a fundamental and essential pillar for successfully constructing [machine learning](#) models that are not only inherently robust but also highly capable of generalization to new, unforeseen data.

Best Practices and Advanced Considerations

While the `createDataPartition()` function significantly simplifies and enhances the process of data splitting, adherence to several best practices and the consideration of additional advanced factors are crucial steps for maximizing the overall effectiveness and scientific rigor of your [predictive](#)

[modeling](#) efforts:

Choosing the Optimal Split Ratio (p): The 80/20 split, as successfully demonstrated, is a widely adopted and excellent standard starting point. However, ratios such as 70/30 or even 90/10 may be more appropriate depending critically on the overall size of your dataset and the inherent complexity of the model you are developing. For exceptionally large datasets, a slightly smaller training percentage might prove sufficient, whereas for smaller datasets, it is often more beneficial to allocate a larger proportion to training or to explore advanced techniques like [k-fold cross-validation](#) to ensure the most robust evaluation possible.

Ensuring Absolute Reproducibility: It is considered an absolute and non-negotiable best practice to consistently invoke `set.seed()` immediately prior to initiating the call to `createDataPartition()`. This specific action guarantees that your random splits are entirely reproducible, which is a critical necessity for effective debugging, for conducting fair and scientifically sound comparisons between different models, and for facilitating the transparent sharing and rigorous verification of your analytical work.

Leveraging Multiple Partitions (`times`): For significantly more rigorous and comprehensively detailed model evaluation, particularly when operating with datasets of limited inherent size or when investigating the stability and consistency of model performance, it is advisable to consider generating multiple, fully independent partitions by explicitly setting the `times` parameter to a value greater than 1. This approach enables [repeated cross-validation](#), a process where a model is trained and tested across several distinct data splits, and its performance metrics are subsequently averaged, leading to a much more reliable and stabilized estimate of generalization error.

Handling Time-Series Data Correctly: For structured [time-series data](#), the application of simple random partitioning methods like `createDataPartition()` is generally inappropriate and can lead to fundamentally misleading results. In this specific context, it is crucial to implement a time-based split, where the training data is derived exclusively from an earlier chronological period, and the testing data is strictly drawn from a subsequent, later period. This methodology accurately simulates real-world prediction scenarios, where models are correctly tasked with forecasting future outcomes based on observations collected in the past.

Mitigating the Risk of Data Leakage: Extreme vigilance is required throughout the process to prevent the occurrence of [data leakage](#). Any data preprocessing steps, which may include operations such as feature scaling, imputation of missing values, or specific types of feature engineering, that rely on statistical information derived from the entire dataset, must be executed *after* the initial data split has been finalized. Alternatively, these crucial operations should be applied separately and entirely independently to the training and testing sets. This prevents any information from the test set from inadvertently influencing the training process, thereby critically preserving the integrity and unbiased nature of the model's final evaluation.

By diligently adhering to these foundational guidelines, you can effectively leverage the robust capabilities of `createDataPartition()` to establish a scientifically sound and exceptionally stable

foundation for all your [machine learning](#) projects within the [R](#) environment. This meticulous and disciplined approach ensures that all your model evaluations are fundamentally fair, statistically accurate, and truly indicative of projected real-world performance.

Additional Resources for R and Machine Learning Proficiency

To further elevate your expertise in data science and [machine learning](#) while utilizing the [R programming language](#), we strongly recommend exploring the following supplementary and authoritative resources:

Consult the official [caret package documentation](#) available on CRAN for the most comprehensive and in-depth information detailing all its functions and full capabilities.

Explore extensive guides and dedicated tutorials focusing on various [data preprocessing](#) techniques in R, which are foundational steps for preparing your data effectively for diverse modeling approaches.

Delve into tutorials that clearly demonstrate the implementation of different [machine learning algorithms](#) in R, building directly upon your foundational understanding of correctly partitioned datasets.

Investigate resources explicitly focused on advanced model evaluation metrics and sophisticated [cross-validation](#) strategies to significantly refine and strengthen your model assessment skills.

These valuable and trusted resources will empower you to substantially expand upon your foundational knowledge of data partitioning and seamlessly integrate it into a complete, highly efficient, and effective [machine learning](#) workflow, ensuring maximum project success.