

How to Easily Replace Missing Values with the dplyr coalesce() Function

Authored by
stats writer

November 27, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Replace Missing Values with the dplyr coalesce() Function*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=100656>

The effective management of data, particularly handling missing values, is a cornerstone of robust statistical analysis and data science. In the environment of the R programming language, the dplyr package provides a powerful suite of functions for data manipulation, and among them, the coalesce() function stands out as an indispensable tool for imputation and data integration.

The primary utility of the coalesce() function is its ability to sequentially scan through a series of inputs (vectors or columns) and return the very first value that is not designated as missing (NA). This functionality is crucial when attempting to reconcile information scattered across multiple variables or when establishing a sensible default value for incomplete observations. If R encounters missingness across all provided arguments at a specific position, the function will simply return NA for that position. This detailed guide explores how to leverage the power of coalesce() through practical, reproducible examples.

Introduction: Understanding Data Coalescing in R

Data cleaning often involves imputing or resolving incomplete records. The concept of "coalescing" refers to the process of merging data points, prioritizing available values over missing indicators. The **coalesce()** function, provided by the essential dplyr package in R, is specifically engineered to handle this operation efficiently. It operates on a principle similar to the SQL standard's COALESCE function, allowing users to define a preference hierarchy among various data sources or columns.

When applying **coalesce()**, the user supplies a sequence of vectors or columns. For every corresponding element position, the function evaluates the inputs from left to right. The moment a non-missing value is encountered, that value is selected and the evaluation stops for that position. This behavior makes it ideal for two common data manipulation scenarios, which we will detail below: replacing NA values within a single vector using a fixed constant, or merging data from alternative columns within a data frame.

These core applications represent highly efficient solutions compared to more verbose conditional statements (like using nested ``ifelse()`` functions). By utilizing **coalesce()**, data cleaning workflows become cleaner, faster, and far more readable.

Syntax and Core Principles of `coalesce()`

The **coalesce()** function accepts any number of arguments, provided they are vectors of the same type or can be coerced to a common type, which is a critical consideration for maintaining data integrity. In its simplest form, you provide the vector requiring correction followed by the value or vector you wish to use as the fill-in.

There are two primary structural patterns for implementing this function:

Method 1: Replace Missing Values in a Vector (Using a Constant Default). This involves providing the target vector first, followed by a single constant value (e.g., 0, "Unknown", or 100) that will be substituted wherever a missing value (NA) is found in the original vector.

library(dplyr)

```
# The vector 'x' is prioritized; if x is NA, use 100.  
coalesce(x, 100)
```

Method 2: Return First Non-Missing Value Across Data Frame Columns (Column Prioritization). This method is used when you have multiple columns, where later columns serve as backups for missing data in earlier columns. This is common when data is redundantly recorded or when prioritizing primary vs. secondary data sources.

library(dplyr)

```
# Check df$A first. If df$A is NA, check df$B next.  
coalesce(df$A, df$B)
```

Method 1: Replacing NA Values in a Single Vector

The simplest and often most common use case for **coalesce()** is replacing all NA entries within a single vector with a predetermined default value. This is a form of simple imputation, essential when downstream functions cannot handle missing data, or when the NA represents a known condition (e.g., zero income, or a score of 100 if the test was missed). This technique is vastly superior to iterating through the vector manually or using complex indexing commands.

When we structure the call as `coalesce(vector_with_NAs, default_value)`, the function effectively performs a check at every index: "Is the value at this position in the first argument available? If yes, keep it. If no (if it is NA), then substitute the value from the second argument." Since the second argument is a constant, it is used whenever a gap is found in the first argument.

Let's demonstrate this practicality by creating a sample vector with several missing values and then applying a fixed imputation value. This example showcases the immediate transformation of incomplete data into a usable, complete vector ready for analysis.

Practical Demonstration: Vector Imputation

The following code illustrates how to use the **coalesce()** function to replace all missing values in a

numeric vector with a specific score of 100. This might be relevant in an academic setting where a missing grade is standardized to a maximum score for certain analyses.

library(dplyr)

```
#create vector of values containing NA (Not Available)
```

```
x <- c(4, NA, 12, NA, 5, 14, 19)
```

```
#replace missing values with 100
```

```
coalesce(x, 100)
```

```
4 100 12 100 5 14 19
```

Observe the resulting output: the initial vector `x` contained two instances of **NA** at the second and fourth positions. These values have been successfully replaced by the constant **100** specified as the secondary argument to **coalesce()**. The other non-missing values (4, 12, 5, 14, 19) remain unchanged, as they satisfied the "first non-missing value" condition. This operation ensures that the resulting vector is complete while retaining all original valid data points.

Method 2: Combining Data Across Multiple Columns

A more complex, yet highly useful, application of **coalesce()** is integrating data from multiple columns within a data frame. Imagine a scenario where a primary measurement (Column A) is often incomplete, but a secondary, less reliable measurement (Column B) is available as a backup. We want to create a new, consolidated column that uses A whenever possible, and falls back to B only when A is missing.

In this structure, the order of arguments is paramount: `coalesce(primary_column, secondary_column)`. The function processes row by row. For each row, it checks the primary column; if it finds a value, that value is taken. If not, it moves to the secondary column and takes its value. This mechanism allows analysts to define explicit data fallback hierarchies, crucial for consolidating information from heterogeneous sources or filling gaps using auxiliary variables.

Suppose we have the following data frame in R, where Columns A and B represent alternate measurements for the same unit:

```
#create data frame
```

```
df <- data.frame(A=c(10, NA, 5, 6, NA, 7, NA),
```

```
B=c(14, 9, NA, 3, NA, 10, 4))
```

```
#view data frame
```

```
df
```

```
A B
```

```
1 10 14
```

```
2 NA 9
```

```
3 5 NA
```

```
4 6 3
```

```
5 NA NA
```

```
6 7 10
```

```
7 NA 4
```

Implementing Column Prioritization with coalesce()

Using the defined data frame `df`, we now apply the **coalesce()** function to create a new Column C, prioritizing values from A over B. This results in an imputed column that leverages the best available data from either source.

library(dplyr)

```
#create new column C that coalesces values from columns A (primary) and B (secondary)
```

```
df$C <- coalesce(df$A, df$B)
```

```
#view updated data frame
```

```
df
```

```
A B C
```

```
1 10 14 10
```

```
2 NA 9 9
```

```
3 5 NA 5
```

```
4 6 3 6
```

```
5 NA NA NA
```

```
6 7 10 7
```

```
7 NA 4 4
```

Analyzing the output reveals the effectiveness of the prioritization:

Row 1: A (10) is non-missing, so C is 10 (even though B is 14).

Row 2: A is NA, so the function falls back to B, resulting in C=9.

Row 3: A (5) is non-missing, so C is 5.

Row 5: Both A and B are NA, meaning no non-missing value was found, hence C remains **NA**.

The resulting column C successfully contains the first non-missing value across columns A and B at each index, demonstrating a highly effective way to merge overlapping data sources while respecting a predefined priority order.

Handling Complete Missingness with a Default Value

As seen in Row 5 of the previous example, if all arguments to **coalesce()** are NA for a specific position, the result remains NA. In many data cleaning contexts, however, we might require a final, ultimate fallback constant to ensure that the resulting vector is entirely free of missing values. This is easily achieved by appending the desired constant as the final argument in the function call.

By adding a third argument, such as the integer 100, we instruct the `coalesce()` function to first check A, then B, and if both are NA, use 100. This guarantees completeness in the output vector, assuming the constant is compatible with the data type of A and B. This multiple-argument capability highlights the flexibility of **coalesce()** as a comprehensive data imputation tool.

library(dplyr)

```
#create new column C that coalesces values from A, B, and finally defaults to 100
df$C <- coalesce(df$A, df$B, 100)
```

```
#view updated data frame
df
```

```
A B C
1 10 14 10
2 NA 9 9
3 5 NA 5
4 6 3 6
5 NA NA 100
6 7 10 7
7 NA 4 4
```

In this updated result, we can clearly see the impact on Row 5. Since both A and B were NA, the function proceeded to the third argument, substituting the NA value with the final default value of **100**. This completes the imputation process, ensuring that the resulting dataset is entirely suitable for analyses that require full data completeness.

Best Practices and Performance Considerations

While **coalesce()** is exceptionally fast and efficient, particularly for large datasets managed within

`dplyr` structures, it is essential to adhere to best practices. Consistency in data types is crucial; attempting to coalesce numeric data with character data will require type coercion, which can sometimes lead to unexpected results if not managed correctly. Always ensure that the vectors being combined are compatible or that the constant default value matches the required output type (e.g., using `coalesce(char_vector, "NONE")` for character data).

Furthermore, while **`coalesce()`** is excellent for simple, sequential imputation, it is not a replacement for complex statistical imputation methods (like mean imputation, K-Nearest Neighbors imputation, or model-based imputation). It excels when the missingness is non-systematic or when a logical, non-statistical fallback value is appropriate, making it a staple function in the initial stages of the data preparation workflow.

Understanding and utilizing the **`coalesce()`** function provides a significant advantage in data wrangling within **R**. Its concise syntax and high performance make it the preferred method for managing missing values through hierarchical prioritization.

Further Reading on Data Manipulation in `dplyr`

The following tutorials explain how to perform other common functions using `dplyr`: