

How to Calculate a Running Total (Cumulative Sum) with dplyr

Authored by
stats writer

December 1, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Calculate a Running Total (Cumulative Sum) with dplyr*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=103288>

The calculation of a cumulative sum is a fundamental operation in data analysis, particularly when tracking totals over time, such as accumulated sales, running balances, or sequential measurements. In the R programming environment, the process is streamlined and highly efficient thanks to the powerful tools provided by the dplyr package. While R provides the base function cumsum(), the integration of this function within the dplyr workflow--leveraging functions like mutate()--allows analysts to calculate running totals easily within a data frame structure without needing complex loops or indices management.

A cumulative sum takes a sequence of numbers and produces a new sequence where each element is the sum of all preceding elements, including the current one. This transformation is invaluable for business intelligence, financial analysis, and scientific research where sequential data processing is required. Using dplyr, we can apply this logic directly to columns, ensuring the calculation remains clean, readable, and highly optimized for speed, which is characteristic of the tidyverse approach to data manipulation. This guide will walk through the two most common scenarios: calculating a running total across an entire dataset and calculating running totals segmented by specific categorical variables (grouping).

The Power of dplyr and the Tidyverse

The dplyr package is the cornerstone of data manipulation in R's tidyverse ecosystem. It offers a standardized and intuitive set of "verbs" designed to operate on tabular data structures efficiently. When dealing with sequential calculations like the cumulative sum, dplyr provides the framework to integrate R's base functions seamlessly into a data pipeline. The elegance of this approach lies in its ability to handle complex operations through a clear sequence of steps, greatly enhancing code maintainability and readability compared to traditional methods.

Central to this workflow is the combination of the mutate() function and the cumsum() function. The mutate() verb is specifically designed to add new variables to an existing data frame, or to transform existing ones, while preserving the existing structure. When cumsum() is applied within mutate(), the calculation is performed element-wise down the specified column, generating the running total column we require. This pairing is crucial for mastering sequential calculations in R.

Furthermore, dplyr heavily relies on the pipeline operator (`%>%`). This operator chains operations together, allowing the output of one function to serve as the input for the next. This structure mimics natural language and thought processes--"take this data, then do this, then do that"--making the resulting code highly intuitive. The following methods demonstrate how this pipeline operator simplifies the calculation of cumulative sums, whether applied globally or within specific groups.

Method 1: Calculating Cumulative Sum of One Column

The simplest application involves calculating the running total across all rows of a specified column. This is often necessary when you are tracking an aggregate metric over time without needing to reset the count based on categories. The general syntax leverages the `mutate()` function alongside the `cumsum()` base R function. The process involves taking your input data structure, piping it into `mutate()`, and defining the new cumulative column within that function call.

The generalized form for this calculation is concise and highly readable. We specify the input data frame (`df`), apply the pipeline operator, and then use `mutate()` to define the new variable, `cum_sum`, by applying `cumsum()` to the target variable, `var1`. This approach ensures that the calculation is executed in sequence, starting from the first row and accumulating the values until the last row is reached. This is an efficient way to convert a series of incremental changes into a running total.

```
df %>% mutate(cum_sum = cumsum(var1))
```

This fundamental technique showcases the power of applying vectorized functions within the `dplyr` context. The base R function `cumsum()` expects a vector as input and returns a vector of the same length containing the running totals. By integrating it into `mutate()`, we guarantee that the calculation aligns perfectly with the rows of the data frame, generating the desired output column without manual row-by-row iteration.

Example 1: Calculating Cumulative Sales Over Time

To illustrate Method 1, consider a scenario where we track daily sales figures for a single entity over an eight-day period. Our goal is to calculate the total accumulated sales as of the end of each day. This requires setting up a basic data frame containing the day number and the sales recorded for that specific day. This setup allows us to clearly observe how the cumulative sum calculation sequentially builds the running total.

We begin by creating our sample data structure in R. Notice how the input data, represented by the `sales` column, provides the incremental value for each observation (day). The cumulative calculation will transform this sequence into a running total, crucial for performance monitoring and forecasting.

```
#create dataset
```

```
df <- data.frame(day=c(1, 2, 3, 4, 5, 6, 7, 8),  
sales=c(7, 12, 10, 9, 9, 11, 18, 23))
```

```
#view dataset
```

```
df
```

```
day sales
1 1 7
2 2 12
3 3 10
4 4 9
5 5 9
6 6 11
7 7 18
8 8 23
```

Upon reviewing the input `data frame`, we see that the first day recorded 7 sales, the second recorded 12, and so on. The cumulative sales column should reflect 7 on Day 1, 19 (7 + 12) on Day 2, 29 (19 + 10) on Day 3, and continue until the final total is achieved on Day 8. Implementing this logic requires loading the `dplyr` library and applying the pipeline syntax we discussed.

The following code block executes the calculation. We are creating a new column named `cum_sales`. This column is defined by applying the base R function `cumsum()` directly to the `sales` column within the context of the `mutate()` verb. This single, concise operation generates the entire sequence of running totals for the dataset.

library(dplyr)

```
#calculate cumulative sum of sales
df %>% mutate(cum_sales = cumsum(sales))
```

```
day sales cum_sales
1 1 7 7
2 2 12 19
3 3 10 29
4 4 9 38
5 5 9 47
6 6 11 58
7 7 18 76
8 8 23 99
```

Interpreting the Results of a Simple Cumulative Sum

The resulting output clearly shows the transformation. The first row of `cum_sales` matches the initial `sales` value (7). Subsequent rows demonstrate the running total: the value 19 in the second row is the sum of the first two daily sales (7 + 12). The final value, 99, represents the total

accumulated sales across all eight days. This calculation is straightforward when dealing with a single continuous sequence, but data often requires more complex partitioning.

When working with large, un-grouped time-series data, the simple application of `cumsum()` is highly efficient. It avoids the performance pitfalls often associated with iterative methods in base R for generating running totals. The speed and clarity offered by `dplyr` make this the preferred method for generating sequential aggregates where grouping is not necessary.

Clarity: The code `df %>% mutate(cum_sales = cumsum(sales))` explicitly states the intent: take the data frame, and calculate a new column containing the running sum of sales.

Efficiency: `cumsum()` is optimized for vector operations in R, ensuring fast computation even for very large datasets.

Integration: The use of `mutate()` ensures the output remains a structured data frame, ready for further analysis or visualization.

Method 2: Calculating Cumulative Sum by Group

In many analytical tasks, the running total must reset whenever a specific categorical variable changes. For example, if you are tracking sales across multiple stores or managing inventory for various product lines, you need the cumulative sum to calculate independently for each group. This requires introducing the `group_by()` verb from the `dplyr` package before applying the cumulative calculation.

The `group_by()` function changes the context of subsequent `dplyr` operations. Once data is grouped, any transformation applied, such as the `mutate()` function containing `cumsum()`, will be executed separately within the confines of each group defined. This powerful feature allows us to avoid writing complex split-apply-combine logic manually, making grouped calculations remarkably straightforward.

The generalized format for calculating a running total by group involves chaining three primary functions using the pipeline operator (`%>%`): first, the data frame; second, `group_by()` specifying the categorical variable(s); and third, `mutate()` applying the `cumsum()` function to the target column. This sequence ensures that the calculation restarts from zero (or the first value of the group) every time a new group begins.

```
df %>% group_by(var1) %>% mutate(cum_sum = cumsum(var2))
```

Example 2: Grouped Cumulative Sales by Store Location

To demonstrate the effectiveness of grouped calculations, let us expand our previous example to include sales data from two different stores, 'A' and 'B', over the same four-day period. The critical requirement here is that the cumulative sales total must reset for Store B, starting its count from Day 1 of its own sequence, independent of Store A's total.

We start by setting up the multi-group data frame. Note that the data is organized sequentially, first listing all observations for Store A, followed by all observations for Store B. The order within the groups is crucial for the cumulative sum calculation, as the running total depends entirely on the existing row order within each partition.

#create dataset

```
df <- data.frame(store=c('A', 'A', 'A', 'A', 'B', 'B', 'B', 'B'),  
day=c(1, 2, 3, 4, 1, 2, 3, 4),  
sales=c(7, 12, 10, 9, 9, 11, 18, 23))
```

#view dataset

```
df
```

```
store day sales
```

```
1 A 1 7
```

```
2 A 2 12
```

```
3 A 3 10
```

```
4 A 4 9
```

```
5 B 1 9
```

```
6 B 2 11
```

```
7 B 3 18
```

```
8 B 4 23
```

In this organized data, we can anticipate the outcome. Store A's cumulative sales should run 7, 19, 29, 38. When the data shifts to Store B (Row 5), the running total must reset. Store B's cumulative sales should start at 9, then become 20 (9 + 11), and so forth. Implementing the grouping logic ensures this reset occurs automatically.

By inserting the `group_by(store)` function into the pipeline operator chain, we instruct dplyr to partition the data logically based on the unique values in the `store` column. The subsequent mutate() call then executes cumsum(sales) independently for Store A's rows and Store B's rows, generating the desired per-group running total.

library(dplyr)

```
#calculate cumulative sum of sales by store
df %>% group_by(store) %>% mutate(cum_sales = cumsum(sales))
```

```
# A tibble: 8 x 4
# Groups: store
store day sales cum_sales
1 A 1 7 7
2 A 2 12 19
3 A 3 10 29
4 A 4 9 38
5 B 1 9 9
6 B 2 11 20
7 B 3 18 38
8 B 4 23 61
```

Verification and Advanced Grouping Concepts

The resulting output successfully demonstrates the grouped cumulative sum. For Store A (Rows 1-4), the `cum_sales` column accumulates to 38. Crucially, in Row 5, where the `store` changes to 'B', the `cum_sales` value resets to 9, which is the sales figure for Store B on Day 1. It then continues to accumulate sales for Store B, reaching a final total of 61. This confirms that the `group_by()` function successfully partitioned the data for the cumulative calculation.

When dealing with time-series or sequential data, it is paramount that the data is correctly ordered within each group before the `cumsum()` function is applied. If the input data frame (`df`) was not already sorted by `day` within each `store`, we would need to include the `arrange()` function in the pipeline operator before `group_by()` and `mutate()`. This ensures the running total accurately reflects the chronological sequence.

For operations involving multiple grouping variables--for instance, grouping by `region` and then by `store` within that region--you simply include all required variables within the `group_by()` call, separated by commas (e.g., `group_by(region, store)`). `dplyr` handles the multi-level partitioning seamlessly, applying the cumulative sum independently for every unique combination of the specified grouping variables. This scalability is why `dplyr` remains the standard for complex data manipulation tasks in R.