

Create a Contingency Table in R

Authored by
Mohammed loot

November 6, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Create a Contingency Table in R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=11326>

A [contingency table](#), frequently known as a cross-tabulation or "crosstab," stands as a cornerstone in quantitative statistical analysis. Its primary purpose is to systematically structure and display the relationship between two or more [categorical variables](#), offering immediate visual insight into their **joint frequencies** and potential associations.

For data scientists and analysts, mastering the analysis of relationships between nominal or ordinal datasets is absolutely crucial. Fortunately, the statistical programming language [R](#) provides highly efficient and powerful native functions specifically designed for generating these tables quickly and accurately. This comprehensive, formal tutorial will meticulously walk you through the entire process: from setting up a realistic dataset to creating, enhancing, and deeply interpreting contingency tables using only core R functions.

The Fundamental Role of Contingency Tables

Contingency tables are indispensable tools because they enable researchers to formally investigate whether two variables are statistically independent or whether they show significant association. The structure is simple yet powerful: each individual cell within the table records the frequency count of observations that simultaneously satisfy a specific combination of categories derived from the two variables under examination.

To illustrate, consider a study on consumer purchasing behavior. A contingency table could cross-reference the product purchased (Variable 1: Television, Radio, Computer) with the region of purchase (Variable 2: Country A, B, C, D). The resulting matrix immediately clarifies patterns, such as identifying which products exhibit the highest sales volume in which specific geographical regions.

While basic frequency tables only summarize the distribution of a single variable, the true analytical strength of the contingency table lies in its capability to expose [joint distributions](#). This foundational capacity is essential for more sophisticated statistical procedures, most notably the [Chi-squared test for independence](#), which is used to formally determine if the observed relationship between the variables is statistically significant and not merely due to random chance. Therefore, understanding the complete structure of these tables--including the rows, columns, and the crucial marginal totals--is the necessary prerequisite for effective data interpretation and rigorous hypothesis testing.

Preparing the Analysis Environment and Sample Data in R

To practically demonstrate the construction of a contingency table, we must first establish a representative sample dataset within the R environment. This dataset is designed to emulate real-world sales records, tracking twenty individual product orders, including the specific product sold and the country where the transaction occurred. We rely on the [data frame](#) structure, which is the standard, most flexible method for organizing and storing tabular data in R.

In the following code block, we use the `data.frame()` function to precisely construct our sales ledger. Note the efficient application of the `rep()` function, which is utilized to repeat category names (TV, Radio, Computer, and Countries A-D) for a specified number of times. This technique populates the columns according to predetermined frequencies, effectively simulating a realistic, non-uniform distribution of sales data that reflects potential real-world biases.

#create data

```
df <- data.frame(order_num = 1:20,  
product=rep(c('TV', 'Radio', 'Computer'), times=c(9, 6, 5)),  
country=rep(c('A', 'B', 'C', 'D'), times=5))
```

```
#view data
```

```
df
```

```
order_num product country
```

```
1 1 TV A
```

```
2 2 TV B
```

```
3 3 TV C
```

```
4 4 TV D
```

```
5 5 TV A
```

```
6 6 TV B
```

```
7 7 TV C
```

```
8 8 TV D
```

```
9 9 TV A
```

```
10 10 Radio B
```

```
11 11 Radio C
```

```
12 12 Radio D
```

```
13 13 Radio A
```

```
14 14 Radio B
```

```
15 15 Radio C
```

```
16 16 Computer D
```

```
17 17 Computer A
```

```
18 18 Computer B
```

```
19 19 Computer C
```

```
20 20 Computer D
```

This resulting data frame, named `df`, now accurately holds 20 distinct observations spanning three variables: a unique order identifier, the purchased product type, and the country of origin for the purchase. Our immediate analytical objective is to summarize the categorical relationship between the `product` and `country` variables.

Generating the Basic Contingency Table using the `table()` Function

The most essential function for generating frequency and contingency tables in [R](#) is the powerful `table()` function. When this function is supplied with two distinct categorical vectors, it automatically computes and tallies the joint occurrences of their unique values, thereby constructing the fundamental contingency structure. This function is remarkably versatile and forms the basis for nearly all categorical data summaries.

The required syntax is straightforward and highly intuitive: we simply pass the two desired variables--in our case, `df$product` and `df$country`--as sequential arguments to the function. It is a standard convention in R that the first argument provided to `table()` typically defines the **row variable**, while the second argument establishes the **column variable** in the final generated output table.

Executing the command below generates the raw frequency count table, which represents the core numerical foundation for all subsequent statistical analysis:

```
#create contingency table  
table <- table(df$product, df$country)
```

```
#view contingency table  
table
```

```
A B C D  
Computer 1 1 1 2  
Radio 1 2 2 1  
TV 3 2 2 2
```

This resulting table meticulously displays the [joint frequency distribution](#). For instance, by locating the intersection of the 'Radio' row and the 'B' column, we find the value 2, which unequivocally indicates that two Radio units were purchased specifically in Country B. While this basic table provides crucial cell frequencies, it lacks the marginal totals necessary for a complete and statistically robust overview.

Enhancing the Table with Marginal Totals using `addmargins()`

For a truly comprehensive statistical understanding, it is absolutely necessary to incorporate the row totals and column totals, which are collectively known as [marginal frequencies](#). These marginal figures provide essential context by illustrating the overall distribution of each variable independently of the other. For example, the row total for 'Computer' tells us the total number of Computers sold across all countries, irrespective of location.

The specialized `addmargins()` function in [R](#) is specifically engineered to calculate and append these crucial sums to any existing table object. The usage is simple and direct: we pass the previously generated `table` object (which contains the raw frequencies) as the sole argument to this function, and R handles the summation automatically.

```
#add margins to contingency table  
table_w_margins <- addmargins(table)
```

```
#view contingency table  
table_w_margins
```

```
A B C D Sum  
Computer 1 1 1 2 5  
Radio 1 2 2 1 6  
TV 3 2 2 2 9  
Sum 5 5 5 5 20
```

The enhanced table, now named `table_w_margins`, is statistically complete. It includes both a 'Sum' row and a 'Sum' column, which dramatically improves both its readability and its analytical utility for any subsequent hypothesis testing. This complete structure facilitates immediate interpretation of both the specific joint frequencies (the inner cells) and the overall marginal distributions (the outer sums).

Interpreting the Results and Drawing Analytical Insights

Effective interpretation of a contingency table complete with margins requires a systematic examination of four distinct, yet interconnected, components: the grand total, the marginal row totals, the marginal column totals, and the individual cell frequencies.

We can systematically interpret the results derived from our sample sales data as follows:

The value located in the bottom right corner (the intersection of the 'Sum' row and 'Sum' column) represents the **grand total**: 20. This value serves as a validation point, confirming the total number of product orders recorded across the entire dataset.

The values running vertically along the right-hand column labeled 'Sum' represent the **row sums**, which detail the [marginal frequency](#) of each product type, irrespective of country:

A total of 5 Computers were ordered worldwide.

A total of 6 Radios were ordered worldwide.

A total of 9 TVs were ordered worldwide.

The values running horizontally along the bottom row labeled 'Sum' represent the **column sums**, which detail the marginal frequency of orders by country, irrespective of product type:

A total of 5 products were ordered from Country A.

A total of 5 products were ordered from Country B.

A total of 5 products were ordered from Country C.

A total of 5 products were ordered from Country D.

The core values situated inside the table represent the **joint frequency**, precisely quantifying the exact number of specific products ordered from each country. For example, the table clearly highlights that 3 TV units were sold in Country A, showcasing a pronounced joint frequency for that particular product-region pairing.

This comprehensive, multi-faceted view allows us to rapidly identify underlying patterns, such as the interesting observation that while the overall sales volume is perfectly evenly distributed across the four countries (5 orders each), the specific product preference within those countries varies quite significantly, pointing toward potential market differences.

Advanced Applications: Converting Frequencies to Proportions

Although raw frequency counts are highly informative, statistical analysts frequently require proportions or percentages to effectively compare distributions, particularly when dealing with samples of differing sizes or when calculating conditional probabilities. [R](#) facilitates this conversion effortlessly by utilizing the `prop.table()` function.

The `prop.table()` function requires two main arguments: the existing contingency table object and an optional margin argument. Specifying `margin = 1` calculates row proportions (percentages relative to the row totals), `margin = 2` calculates column proportions (percentages relative to the column totals), and omitting the margin argument calculates proportions relative to the overall grand total.

For instance, if our objective is to determine what proportion of all sales within each country consisted specifically of Televisions, we would calculate the column proportions, standardizing the counts relative to the country totals:

```
# Calculate column proportions (percentages relative to the country total)
```

```
prop_col <- prop.table(table, margin = 2)
```

```
# view the proportional table
```

```
prop_col
```

```
A B C D
```

```
Computer 0.20 0.20 0.20 0.40  
Radio 0.20 0.40 0.40 0.20  
TV 0.60 0.40 0.40 0.40
```

In examining this proportional table, we can clearly see that 60% of all sales recorded in Country A were TVs (derived from $3/5 = 0.60$). Conversely, in Country D, Computer sales represented 40% of the total ($2/5 = 0.40$). Analyzing data in terms of relative proportions provides crucial context that extends far beyond simple raw counts, making this step absolutely critical in any robust comparative analysis.

Summary of Essential R Functions for Contingency Analysis

Successfully mastering the creation and interpretation of contingency tables in [R](#) hinges on a solid understanding of the synergistic interplay between three primary functions, each of which serves a distinct and vital purpose in the data aggregation pipeline:

`table(var1, var2)`: This serves as the initial, core function for generating the raw frequency counts, strictly defining the rows and columns based on the input categorical variables. It is the mandatory starting point for all cross-tabulation analysis.

`addmargins(table_object)`: This function is utilized immediately after `table()` to accurately calculate and append the essential [marginal sums](#) (both row totals and column totals) and the grand total to the table, ensuring comprehensive frequency context.

`prop.table(table_object, margin)`: This function converts the absolute frequency counts into relative proportions. This conversion capability is vital for comparing relative frequencies, addressing questions of conditional probability, and understanding percentage distributions across different categories.

These three specialized tools collectively empower analysts to conduct robust statistical exploration of the complex relationships between [categorical variables](#), efficiently transforming raw data into reliable, actionable insights for informed decision-making.