

# How to Filter Columns in dplyr Based on a Starting String

Authored by  
**stats writer**

January 16, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Filter Columns in dplyr Based on a Starting String*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=126299>

Filtering data based on specific textual patterns is a fundamental task in data analysis. When working with R, the powerful combination of the dplyr package for data manipulation and the stringr package for string operations provides an elegant solution for implementing a "starts with" filter. This method allows you to select specific rows from a data frame where the values in a designated column begin with a predefined prefix or sequence of characters. Although dplyr offers `starts_with()` as a helper function primarily for selecting columns by name, achieving prefix matching within the rows of a column requires utilizing the `filter()` verb in conjunction with specialized string detection functions.

This technique is indispensable when dealing with large datasets where categories or descriptions share common prefixes, such as identifying all product codes starting with "A-", or in our case, identifying player positions that begin with "backup" versus "starting". The synergy between dplyr's expressive syntax and stringr's robust text handling capabilities makes complex filtering tasks straightforward and highly readable. Mastering this specific filtering mechanism ensures that your data preparation is precise and efficient, paving the way for deeper statistical analysis.

To implement this logic effectively, we leverage the concept of regular expressions (regex), which are essential for pattern matching in text. Specifically, the caret symbol (^) acts as an anchor in regex, signifying the absolute start of a string. By incorporating this simple yet powerful symbol into our filtering condition, we guarantee that only rows where the target column entry begins with the desired pattern are included in the resulting output data frame. The following sections will detail the necessary packages, the crucial syntax elements, and provide practical examples to illustrate this powerful filtering method in R.

## The Essential Combination: dplyr and stringr

To successfully implement prefix-based row filtering in R, we require functionality from two core packages: dplyr, which provides the highly efficient data pipeline structure (using the pipe operator `%>%`) and the `filter()` function, and stringr, which specializes in user-friendly string manipulation functions, including pattern detection. The stringr function `str_detect()` is the key component here, as it tests whether a string contains a specific pattern and returns a logical vector (TRUE/FALSE) that `filter()` can then use to select rows.

The standard syntax involves piping your target data frame into the `filter()` function. Within `filter()`, you call `str_detect()`, specifying the column you wish to inspect (e.g., `position`) and the regular expression pattern to search for. The pattern must include the caret symbol (^) immediately preceding the desired prefix. This combination tells stringr to look for the pattern only at the absolute start of the string value in that specific column.

The following structural example demonstrates how to set up this operation. Note that we must first load both libraries to ensure the functions are accessible within the R session. This particular

command sequence is designed to filter the data frame named **df**, keeping only those rows where the value in the **position** column begins exactly with the string "back":

You can use the following basic syntax to filter for rows where a column starts with a certain pattern:

```
library(dplyr)
```

```
library(stringr)
```

```
df %>%  
filter(str_detect(position, "^back"))
```

This particular example filters the data frame named **df** to only show the rows where the **position** column starts with the string "back."

**Note on Regular Expressions:** In regex, the **^** symbol indicates the beginning of a string. It is this anchoring feature that transforms a general string search into a highly specific "starts with" filter. Without the caret, `str_detect("backup", "back")` would be TRUE, but it would also be TRUE for a string like "halfback", which might not be the desired behavior if strict prefix matching is required.

## Creating the Sample Data Frame in R

To demonstrate the utility and precision of the `str_detect()` filtering method, we will first construct a simple sample dataset in R. This data frame, named **df**, represents information about six fictional basketball players, specifying their identifiers and their specific positions. Notice that the positions clearly fall into two categories based on their prefix: "starting" and "backup". This structure makes it an ideal scenario for applying our "starts with" filter.

The **position** column contains variations like `starting_guard`, `backup_center`, and `starting_forward`. Our primary objective is to use a filter based on the string prefix to quickly isolate all players who are designated as "backup" players, excluding all those designated as "starting" players. This requires us to focus our pattern matching exclusively on the beginning of the string.

Here is the R code used to generate and inspect our sample dataset:

```
#create data frame  
df <- data.frame(player=c('A', 'B', 'C', 'D', 'E', 'F'),  
position=c('starting_guard', 'starting_center', 'backup_guard',  
'backup_center', 'starting_forward', 'backup_forward'))
```

```
#view data frame
df

player position
1 A starting_guard
2 B starting_center
3 C backup_guard
4 D backup_center
5 E starting_forward
6 F backup_forward
```

As observed in the output above, the data frame **df** is correctly populated, containing six rows and two columns. We now have a robust structure upon which we can apply the `dplyr` and `stringr` filtering combination to isolate subsets based on their position prefixes.

### Case Study 1: Filtering Rows Starting with a Specific Word Prefix ("back")

Our first filtering objective is to extract all players whose position starts with the word "back". This means we are targeting rows associated with Player C, D, and F. The critical element in our pattern is `^back`. The caret (^) ensures that the match must begin at the very first character of the string, guaranteeing that "backup\_guard" matches, while a hypothetical entry like "fullback\_team" (if present) would also match, but something containing "back" internally, like "midfield\_back", would be correctly excluded.

We execute the filtering operation by loading the necessary libraries and utilizing the pipeline operator (`%>%`). The input data frame **df** is channeled into `filter()`, where the `str_detect()` function evaluates the **position** column against our predefined regex pattern. Only rows returning TRUE from the `str_detect` call are retained in the final output.

This method offers superior efficiency and clarity compared to traditional base R subsetting methods involving functions like `substr()` or `grep()`, especially within a complex data manipulation pipeline powered by `dplyr`. The code snippet below illustrates the precise implementation and the resulting filtered data frame:

```
library(dplyr)
library(stringr)
```

```
#filter data frame to only contain rows where position column starts with "back"
df %>%
filter(str_detect(position, "^back"))
```

```
player position
1 C backup_guard
2 D backup_center
3 F backup_forward
```

Upon reviewing the output, we confirm that the resulting data frame successfully isolates the three rows where the **position** string starts with the prefix "back." This result validates the effective use of the `^` anchor within `str_detect()` to achieve highly targeted prefix filtering.

## Case Study 2: Filtering Rows Based on a Single Initial Character ("s")

The power of the `^` anchor is not limited to searching for lengthy prefixes; it is equally effective for filtering based on a single starting character. In our basketball player dataset, suppose we wish to identify all players whose position begins with the letter "s", which corresponds to the "starting" positions. This is a common requirement when standardizing data entries or extracting information based on initial categorization markers.

To achieve this, we simply modify the `regex` pattern within `str_detect()` to `^s`. This pattern strictly checks the very first character of the string in the **position** column. Because "starting\_guard", "starting\_center", and "starting\_forward" all begin with 's', they will all be captured by this precise filter, while the "backup" positions (starting with 'b') will be excluded.

This demonstrates the flexibility of using `stringr` functions within the `dplyr` framework. Whether the desired prefix is a lengthy word or a solitary character, the underlying mechanism remains robust and consistent: the `filter()` function executes the logical test provided by `str_detect()` using the string-start anchor. The following code provides the implementation and result for this single-character filtering scenario:

```
library(dplyr)
library(stringr)

#filter data frame to only contain rows where position column starts with "s"
df %>%
  filter(str_detect(position, "^s"))

player position
1 A starting_guard
2 B starting_center
3 E starting_forward
```

We can see that the resulting data frame only contains rows where the string in the **position** column starts with the letter **s**. This capability ensures analysts can quickly partition data based on initial characters, regardless of string complexity.

## Handling Case Sensitivity and Alternative Methods

It is important to understand that the standard implementation of `str_detect()` is case-sensitive. If our sample data frame contained "Starting\_guard" (with a capital S) alongside "starting\_guard" (lowercase s), the pattern `^s` would only match the lowercase entries. To perform a case-insensitive "starts with" filter, you must adjust the function call within stringr.

The `str_detect()` function includes an optional argument, `ignore.case`, which can be set to `TRUE` to enable case-insensitive matching. Alternatively, you could use `stringr::str_to_lower()` or `stringr::str_to_upper()` to normalize the entire column before filtering. However, setting the `ignore.case` argument is generally the cleanest method for this specific purpose, allowing the pattern `^s` to match both "Starting" and "starting".

For example, to filter for strings starting with "start", regardless of capitalization, the syntax would be: `filter(str_detect(position, "^start", ignore.case = TRUE))`. This flexibility allows for robust filtering even when data entry standards are inconsistent regarding capitalization. Failure to account for case sensitivity is a common source of data filtering errors in analytical workflows.

## Alternative Approach: Using Base R or String Manipulation

While the `dplyr/stringr` combination is highly recommended for its readability and integration into the tidyverse ecosystem, it is worth noting that the "starts with" filtering can also be achieved using base R functions or alternative methods. One common base R approach involves using the `grepl()` function (Global Regular Expression Print Logical) in conjunction with the subset operator

Another alternative, particularly useful when dealing with fixed-width prefixes and avoiding the complexity of regex, is using the `substr()` function. If you know the exact length of the prefix (e.g., 4 characters for "back"), you can extract the first four characters and compare them directly:

```
# Base R equivalent using substr()  
df
```

While these base R methods are functional, they often interrupt the clean data flow established by the dplyr pipeline. The `str_detect()` approach within `filter()` provides a more cohesive and generally more readable solution for most modern R analysis projects.

## Summary of the "Starts With" Filtering Technique

The effective implementation of a "starts with" filter in R is best achieved by integrating the string manipulation capabilities of the `stringr` package, specifically `str_detect()`, with the data manipulation infrastructure of `dplyr`'s `filter()`. The fundamental key to this operation is the use of the regular expression anchor `^` (caret), which locks the pattern search to the beginning of the string value. This ensures absolute precision in identifying rows based purely on the prefix.

By using the pipeline operator (`%>%`), data analysts can seamlessly incorporate this filtering step into complex transformations, ensuring data subsets are correctly identified and isolated before further aggregation or visualization. This consistency and power are primary reasons why this methodology is preferred in the modern R ecosystem.

To summarize the steps for filtering a data frame `df` based on a prefix `P` in column `C`:

Load both the `dplyr` and `stringr` libraries.

Use the pipe operator: `df %>%`.

Call the filter function: `filter()`.

Inside filter, call `str_detect`: `str_detect(C, "^P")`.

## Related Functions and Next Steps in dplyr

Once you have mastered the technique for filtering rows based on prefixes, you may find related operations useful for comprehensive data cleaning and preparation. For instance, if you needed to filter for strings that **end with** a certain suffix, you would use the dollar sign anchor (`$`) in your regex pattern (e.g., `"guard$"`). Similarly, if you wanted to select columns whose names start with a specific prefix, you would use the `starts_with()` helper function directly within `select()` or `rename_with()`.

The following common functions in `dplyr` represent essential steps in the data workflow:

`select()`: Used for keeping or discarding specific columns.

`mutate()`: Used for adding new columns or transforming existing ones.

`group_by()` and `summarize()`: Used together to calculate summary statistics for defined groups within the data.

`arrange()`: Used for sorting the data frame based on column values.

By combining these powerful verbs with precise string filtering techniques, you gain complete control over your data manipulation tasks in R.