

# How to Filter PySpark DataFrames for Values Not Containing a Substring

Authored by  
**stats writer**

February 10, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Filter PySpark DataFrames for Values Not Containing a Substring*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=129975>

## Introduction to Advanced Data Filtering in PySpark

In the contemporary landscape of **big data** analytics, PySpark has emerged as an indispensable framework for managing and processing vast datasets across distributed computing environments. As the Python API for Apache Spark, it combines the simplicity and versatility of the Python programming language with the robust, high-performance capabilities of the Spark engine. One of the most fundamental yet critical operations in any data pipeline is the ability to selectively isolate or exclude specific records based on complex criteria. Mastering the nuances of the DataFrame.filter method is essential for any data engineer or scientist looking to refine their datasets for downstream analysis or **machine learning** applications.

The process of data cleaning often requires more than just identifying what information to keep; it frequently necessitates the identification and removal of "noise" or irrelevant records that match specific, unwanted patterns. While positive filtering--selecting rows that contain a specific substring--is intuitive, the inverse operation is equally vital. PySpark provides a sophisticated set of tools to perform "not contains" operations, allowing users to efficiently prune their DataFrames. This capability is particularly useful when dealing with messy, real-world data where certain keywords might signify errors, test entries, or categories that fall outside the scope of a particular study.

By leveraging logical operators and built-in column functions, PySpark users can construct highly readable and performant queries to handle substring exclusion. This guide provides an in-depth exploration of how to implement these negative filters, focusing on the use of the logical NOT operator in conjunction with the **contains** method. We will examine the underlying mechanics of these operations, their syntax, and provide a comprehensive walk-through of a practical example to ensure you can apply these techniques to your own large-scale data processing tasks with confidence and precision.

### The Mechanics of the Negative Filter and the Tilde Operator

At the core of PySpark's filtering logic is the ability to evaluate Boolean expressions for every row in a distributed dataset. When we wish to exclude rows based on a substring, we essentially need to flip the result of a standard inclusion test. In PySpark, this inversion is most commonly achieved using the **tilde (~)** symbol, which serves as the logical negation operator. When applied to a Column expression that returns a Boolean value, the tilde converts all **true** results to **false** and all **false** results to **true**, effectively allowing the **filter** function to discard the rows that would have otherwise been selected.

The **contains()** method itself is a member of the Column class. It checks whether a specified string literal is present within the values of the target column. While this method is highly efficient for simple substring matching, it is important to remember that it operates on a per-character basis

and is strictly **case-sensitive**. Understanding this behavior is crucial because a mismatch in casing can lead to rows being inadvertently retained in the dataset, potentially skewing the results of subsequent data analysis or statistical modeling.

Furthermore, the use of the tilde operator is preferred over other negation methods in PySpark due to its conciseness and compatibility with Spark's query optimizer, known as Catalyst. When you chain operations or use logical operators, the Catalyst Optimizer generates an efficient execution plan that minimizes data shuffling across the cluster. By utilizing these native Spark functions, you ensure that your code remains performant even when processing billions of records, maintaining the scalability that makes PySpark so powerful in enterprise environments.

## PySpark: Filter for "Not Contains"

**To implement a "Not Contains" logic within your data pipeline, you can utilize the tilde operator as a prefix to the contains method. This syntax is highly expressive and integrates seamlessly into the standard PySpark DataFrame API:**

```
#filter DataFrame where team does not contain 'avs'  
df.filter(~df.team.contains('avs')).show()
```

**The practical application of this syntax is demonstrated in the following detailed example, which illustrates the transformation of a raw dataset into a refined subset.**

**Example: Implementing a "Not Contains" Filter in PySpark**

**To begin our demonstration, we must first establish a SparkSession, which serves as the entry point for any**

**PySpark application. In this scenario, we will construct a synthetic dataset representing performance metrics for various professional basketball teams. This dataset includes a categorical 'team' column and a numerical 'points' column, providing a clear structure for testing our filtering logic:**

```
from pyspark.sql import SparkSession  
spark = SparkSession.builder.getOrCreate()
```

```
#define data
```

```
data = ,
```

```
,
```

```
,
```

```
,
```

```
,
```

```
,
```

```
,
```

```
,]
```

```
#define column names
```

```
columns =
```

```
#create dataframe using data and column names
```

```
df = spark.createDataFrame(data, columns)
```

```
#view dataframe
```

```
df.show()
```

```
+-----+-----+  
| team|points|  
+-----+-----+  
| Mavs| 14|  
| Nets| 22|  
| Nets| 31|  
| Cavs| 27|  
| Kings| 26|  
| Spurs| 40|  
|Lakers| 23|  
| Spurs| 17|  
+-----+-----+
```

Executing the Substring Exclusion Logic

With the DataFrame successfully initialized, we can now apply the negative filter. Suppose our objective is to exclude any team whose name includes the substring "avs". This would typically include teams like the Mavericks and the Cavaliers. By applying the tilde operator to the column's contains method, we can generate a new view of the data that only includes

teams that do not match this specific pattern:

```
#filter DataFrame where team does not contain 'avs'  
df.filter(~df.team.contains('avs')).show()
```

```
+-----+-----+  
| team|points|  
+-----+-----+  
| Nets| 22|  
| Nets| 31|  
| Kings| 26|  
| Spurs| 40|  
|Lakers| 23|  
| Spurs| 17|  
+-----+-----+
```

Upon reviewing the resulting output, it is evident that the logic has been applied correctly. The records for 'Mavs' and 'Cavs' have been successfully omitted from the display. This demonstrates the precision of the filter operation when combined with logical negation, ensuring that only the relevant data points--those not containing the specified string--are passed through to the next stage of the data pipeline.

## Analysis of the Filtered Results

The effectiveness of this operation lies in its ability to identify substrings regardless of their position within the string. In our example, the substring "avs" appeared at the end of both "Mavs" and "Cavs". Because the `contains` function searches the entire length of the string, it successfully flagged these records for exclusion. This is a powerful feature for data cleaning, as it allows for the removal of wide swaths of data that may share a common naming convention or error code without requiring an exact match for the entire string.

It is also important to observe that other teams, such as the 'Nets' or 'Spurs', remained entirely unaffected by the filter. This level of granularity ensures that your data manipulation remains targeted. When working with large schemas containing hundreds of columns and millions of rows, such precision is necessary to maintain the integrity of your findings. By isolating these specific patterns, you can focus your computational resources on the data that truly matters for your business logic or scientific inquiry.

## Addressing Case Sensitivity in PySpark String Operations

A critical technical detail to keep in mind is that the `contains` function is inherently case-sensitive. In the context of our basketball example, if the filter had been defined to look for the uppercase substring "AVS", the results would have been significantly different. Because "Mavs" and "Cavs" utilize lowercase letters for that specific sequence, a case-sensitive search for "AVS" would return false for those rows, and the negation would therefore keep them in the final DataFrame.

To mitigate issues arising from inconsistent casing, data practitioners often employ data normalization techniques. One common approach is to transform the entire column to lowercase or uppercase using the `lower()` or `upper()` functions before performing the comparison. This ensures that the filter is robust and accounts for variations in how data might have been entered into the source system. For instance, `df.filter(~lower(df.team).contains('avs'))` would provide a more foolproof way to ensure all variations of the substring are captured and excluded.

Advanced Pattern Matching with RLIKE and LIKE

While the `contains` method is excellent for static

substrings, there are scenarios that require more complex matching logic. For these instances, **PySpark** offers the `like` and `rlike` functions. The `like` function allows for the use of SQL-style wildcards, such as the percent sign (%) to represent zero or more characters. This is useful if you want to exclude strings that start or end with a certain pattern but might contain other characters in between.

The `rlike` function takes this a step further by supporting regular expressions (Regex). Regular expressions allow for incredibly sophisticated pattern matching, such as identifying strings that contain digits, specific sequences of vowels, or complex structural formats like email addresses or phone numbers. By applying the tilde operator to an `rlike` expression, such as `df.filter(~df.team.rlike('^M|s$'))`, you can exclude rows based on highly specific architectural rules, making it a cornerstone technique for advanced data engineering.

#### Performance Considerations for Large Scale Filtering

When executing filters on massive datasets, it is important to consider the underlying architecture of

**Apache Spark**. Filtering is a "narrow transformation," meaning it does not require data to be shuffled across different nodes in the cluster. Each partition of the data can be filtered independently on the executor where it resides, which makes this operation highly efficient. However, the complexity of the string matching logic can still impact the overall execution time.

Using `contains` is generally faster than using `rlike` because regular expression engines require more computational overhead to parse and execute the matching patterns. Therefore, it is considered a best practice to use the simplest tool that satisfies your requirements. If a basic substring match is sufficient, `contains` is the optimal choice. Additionally, placing filters as early as possible in your **Directed Acyclic Graph (DAG)** can significantly improve performance by reducing the volume of data that subsequent transformations--especially wide transformations like joins or group-bys--must process.

Summary of Best Practices for Substring Exclusion

In conclusion, the ability to filter for values that do not contain a specific substring is a powerful feature in

**PySpark** that facilitates cleaner, more relevant data sets. By utilizing the tilde (~) operator in conjunction with the contains method, users can implement logical negation that is both readable and highly efficient. Whether you are performing exploratory data analysis or building complex automated data pipelines, understanding these filtering techniques is essential for maintaining high standards of data quality and operational efficiency.

As you progress in your journey with PySpark, remember to account for case sensitivity and consider using normalization functions like lower() to make your filters more resilient. For more complex requirements, do not hesitate to explore rlike and regular expressions, but always keep performance in mind. By following the structured approach outlined in this guide, you can ensure that your data processing remains accurate, scalable, and professional.

The following tutorials and documentation resources provide further insights into performing other common tasks and optimizations within the PySpark ecosystem, helping you to further expand your technical expertise

in big data management.

ARABPSYCHOLOGY.COM