

Learning Pandas: Filtering DataFrames – Selecting Rows Based on Value Ranges

Authored by
Mohammed Iooti

November 15, 2025

RECOMMENDED CITATION

Mohammed Iooti (2025). *Learning Pandas: Filtering DataFrames – Selecting Rows Based on Value Ranges*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=2527>

In the demanding field of [data analysis](#) and high-volume data manipulation, one task remains perpetually fundamental: efficiently filtering datasets to isolate specific, meaningful subsets of information. When working with [tabular data](#) using [Pandas](#), the cornerstone Python library for data science, it is frequently necessary to select rows where a value in a designated column falls precisely within a defined numerical or chronological range. This operation is indispensable for a variety of analytical objectives, including the identification of anomalies or outliers, the segmentation of populations for targeted studies, and the preparation of clean, targeted data subsets necessary for machine learning model training.

Historically, achieving this range filtering manually involved chaining multiple conditional comparison operators, such as `(df >= lower_bound) & (df <= upper_bound)`. While functional, this approach quickly becomes verbose, cumbersome, and error-prone, especially when dealing with complex datasets or when the boundary requirements shift between inclusive and exclusive thresholds. Recognizing this common bottleneck, Pandas provides an elegant, highly optimized, and exceptionally readable solution tailored specifically for this purpose: the [.between\(\) method](#). This function drastically simplifies the process of [DataFrame](#) filtering, allowing data professionals to intuitively specify both the lower and upper constraints for inclusion, streamlining complex data preparation workflows.

The core syntax for employing the `.between()` method utilizes Pandas' powerful [Boolean indexing](#) capabilities. The method is applied directly to the target Series (column), which subsequently generates a Boolean mask used to filter the parent DataFrame. This structure ensures concise and effective code:

```
df_filtered = df.between(25, 35)]
```

The example above creates the `df_filtered` result, which contains all rows from the original `df` where the value in the `'points'` column falls inclusively between 25 and 35. Furthermore, if the analytical goal is to isolate data points that are **outside** this defined range--a common approach for identifying statistical outliers or values that breach established thresholds--Pandas allows for an immediate inversion of the selection. This inversion is easily executed by prefixing the entire Boolean condition with the bitwise [tilde](#) (`~`) operator:

```
df_filtered = df.between(25, 35)]
```

The subsequent sections will provide a deep dive into the mandatory and optional parameters of the [.between\(\) method](#), illustrate its implementation through practical, real-world coding examples, and detail the necessary best practices for achieving precision and efficiency in data filtering across diverse analytical scenarios.

Understanding the `.between()` Method Parameters

The `.between()` method is fundamentally designed as a Series method within [Pandas](#), meaning its operation is concentrated on a single dimension of data--a column--nested within a larger [DataFrame](#). When applied to a column, its singular purpose is to evaluate, row by row, whether the value meets the defined range criteria. The resulting output of this evaluation is not the filtered data itself, but rather a [Boolean Series](#) (often referred to as a mask), where a value of `True` signifies that the corresponding row's data falls within the stipulated range, and `False` indicates that it does not. This mask is the critical component used for Boolean indexing, which filters the original DataFrame to yield only the rows marked `True`.

Grasping the method signature is vital for leveraging its full analytical power. The standard structure is defined as `Series.between(left, right, inclusive='both')`. The three core parameters--`left`, `right`, and `inclusive`--govern the entirety of the filtering process and must be handled meticulously to satisfy precise analytical requirements. The first two parameters, `left` and `right`, establish the boundaries of the selection interval. Specifically, the `left` parameter invariably defines the **lower bound** of the range, meaning values in the Series must be greater than this value (or equal to it, depending on the third parameter). Conversely, the `right` parameter establishes the **upper bound**, requiring Series values to be less than this value (or equal to it).

The most influential and critical aspect of the `.between()` method is the optional but powerful `inclusive` parameter, which explicitly dictates how the boundaries themselves are handled during the comparison process. By default, `inclusive` is set to the string `'both'`, which treats the operation as inclusive of both the lower and upper limits. However, advanced [data analysis](#) frequently demands stricter or more fluid boundary conditions. The `inclusive` parameter accepts four distinct string values, providing granular control over the selection logic, enabling the definition of open, closed, or semi-open intervals:

'both' (Default): This setting ensures that both the `left` and `right` bounds are included in the final selection. Mathematically, this corresponds to the closed interval condition: `left <= value <= right`.

'neither': This setting strictly excludes both the `left` and `right` bounds, necessitating that the value must fall strictly between the limits. The mathematical condition is: `left < value < right`.

'left': This setting includes the `left` bound but strictly excludes the `right` bound. This defines a semi-open interval where the condition is: `left <= value < right`.

'right': This setting strictly excludes the `left` bound but ensures the inclusion of the `right` bound. The condition is: `left < value <= right`.

Mastering the nuanced application of the `inclusive` parameter is essential for guaranteeing precise and accurate data filtering, ensuring that the resulting subset of data perfectly aligns with the required analytical scope.

Practical Example: Filtering Rows with Default Inclusion

To clearly demonstrate the straightforward application of the [.between\(\) method](#), we will examine a scenario based on processing sports statistics. Imagine we are tasked with analyzing player performance data in a professional basketball league, and our goal is to isolate players who achieved a moderate scoring output--those who are neither the absolute top performers nor those who scored minimally. This requirement necessitates filtering the dataset based on a specific score range, inclusive of the boundaries.

We begin this process by initializing a sample [DataFrame](#). This DataFrame, named `df`, is structured with two key columns: `'team'`, which registers the team affiliation, and `'points'`, which contains the score recorded by each player in a hypothetical game. This setup provides a solid foundation for testing our filtering criteria:

```
import pandas as pd
```

```
#create DataFrame
```

```
df = pd.DataFrame({'team': ,  
'points': })
```

```
#view DataFrame
```

```
print(df)
```

```
team points
```

```
0 Mavs 22
```

```
1 Mavs 28
```

```
2 Nets 35
```

```
3 Nets 34
```

```
4 Heat 29
```

```
5 Heat 28
```

```
6 Kings 23
```

With the DataFrame successfully initialized, we proceed to apply the `.between()` method. Our specific requirement is to select all rows where the value in the `'points'` column is inclusively between 25 and 35. Since the default value for the `inclusive` parameter is `'both'`, we only need to supply the lower bound (25) and the upper bound (35) to the method, allowing the Pandas engine to manage the inclusive comparison automatically and efficiently.

```
#select rows where value in points column is inclusively between 25 and 35
```

```
df_filtered = df.between(25, 35)]
```

```
#view filtered DataFrame
```

```
print(df_filtered)
```

```
team points
```

```
1 Mavs 28
```

```
2 Nets 35
```

```
3 Nets 34
```

```
4 Heat 29
```

```
5 Heat 28
```

The resulting `df_filtered` DataFrame clearly illustrates the effect of the default inclusive behavior. Rows corresponding to scores of 28, 29, 34, and 35 are all successfully included in the output. Crucially, the boundary value of 35 is retained, confirming the method's behavior, while scores such as 22 and 23 are correctly excluded because they fall below the specified lower bound of 25. This successful outcome validates that the default application of `.between()` accurately filters the data to the specified closed interval.

Refining Selection: Achieving Strict Exclusion

While the default `inclusive='both'` setting serves most general filtering tasks, sophisticated [data analysis](#) often mandates strict boundary exclusion. The true utility of the `inclusive` parameter lies in its capacity to provide this granular control, allowing analysts to tailor selections with absolute precision. For example, if our objective shifts to identifying players whose scores are strictly greater than the minimum threshold (25) but strictly less than the maximum threshold (35), we must explicitly exclude both boundaries.

To achieve a selection that is **strictly between** 25 and 35, meaning scores of exactly 25 and 35 must be discarded, we must modify our previous code snippet by setting the parameter `inclusive='neither'`. This instruction tells [Pandas](#) to select only those values that are strictly greater than the left bound AND strictly less than the right bound. This adjustment is essential when defining true open intervals, a frequent requirement in rigorous statistical or performance metric evaluations.

```
#select rows where value in points column is strictly between 25 and 35
```

```
df_filtered_exclusive = df.between(25, 35, inclusive='neither')]
```

```
#view filtered DataFrame
```

```
print(df_filtered_exclusive)
```

```
team points
1 Mavs 28
3 Nets 34
4 Heat 29
5 Heat 28
```

By carefully examining the output of `df_filtered_exclusive`, the effect of the parameter change is immediately evident: the player who scored exactly 35 points (original index 2) is now excluded from the results. This outcome stands in clear contrast to the previous result where `inclusive='both'` was used. This example powerfully highlights the flexibility provided by the `inclusive` parameter, which empowers analysts to define open, closed, or semi-open intervals (by using `'left'` or `'right'`) precisely as required by the specific data processing task.

Filtering Rows Outside a Specified Range for Anomaly Detection

In many high-stakes, real-world analytical scenarios, the primary focus shifts away from identifying data points that fall *within* a standard operating range, towards isolating those that fall *outside* it. This technique is invaluable for applications such as quality control, robust anomaly detection, and the identification of extreme performance metrics (both high and low). [Pandas](#) facilitates this critical inversion of selection criteria through the use of the bitwise NOT operator, symbolized by the [tilde](#) (`~`).

As previously established, the `.between()` method generates a [Boolean Series](#) where `True` corresponds to values located inside the specified range. When the tilde operator (`~`) is applied to this resulting Series, it effectively flips every logical value: `True` becomes `False`, and `False` becomes `True`. Consequently, applying this inverted Boolean mask to the [DataFrame](#) selects all rows that originally failed the `.between()` condition, thereby identifying "everything else"--the outliers--outside the defined 25 to 35 range.

#select rows where value in points column is not between 25 and 35 (inclusive)

```
df_filtered = df.between(25, 35)]
```

```
#view filtered DataFrame
print(df_filtered)
```

```
team points
0 Mavs 22
6 Kings 23
```

As clearly demonstrated by the output, only the players scoring 22 and 23 points are included in

the resulting DataFrame. These are the two data points that fell below the lower bound of 25. Had the original dataset contained scores greater than 35, those records would also have been included in this result. This technique offers a concise and highly readable alternative to manually constructing complex chained logical expressions (e.g., `(df < 25) | (df > 35)`), making the analytical intent--filtering for outliers--immediately obvious to any stakeholder reading the code.

Important Considerations and Best Practices

While the `.between()` method is exceptionally robust, maximizing its efficiency and ensuring reliable results mandates adherence to certain best practices, particularly regarding data types, performance, and the explicit handling of edge cases like missing values. Neglecting these critical details can invariably lead to unexpected filtering results or runtime errors in production environments.

Data Types Compatibility: The primary strength of `.between()` is rooted in its application to continuous data. It is specifically optimized for [numeric data](#) types (such as integers, floats, and decimals) and chronological data represented by `datetime` objects. When filtering temporal data, it is absolutely essential to confirm that the target column has been correctly converted to the appropriate Pandas `datetime64` type. If the method is mistakenly applied to string data, `.between()` will execute a lexicographical (alphabetical) comparison, which rarely aligns with the intuitive concept of a numerical range, potentially resulting in misleading or incorrect analytical conclusions if not managed with vigilance.

Handling Missing Values (NaN): A critical operational detail of `.between()` is its default behavior concerning [missing values \(NaN\)](#). By inherent design, any row containing a `NaN` in the column being evaluated for the range will automatically yield a `False` value in the resulting Boolean mask. This means that rows with missing data are systematically excluded from the filtered results. Analysts must be keenly aware of this behavior. If your analysis requires that rows with missing values be included or treated in a specific manner (e.g., if you are explicitly searching for missing data), you must first preprocess your data--perhaps by imputing `NaNs` with a proxy value or by creating a separate conditional filter--before applying `.between()`.

Performance Implications: For [data science](#) professionals dealing with extremely large [DataFrame](#) objects, optimization is paramount. Utilizing `.between()` is universally considered the recommended approach for range filtering because its underlying implementation is highly optimized, leveraging efficient C code. This method consistently outperforms the explicit chaining of comparison operators (e.g., `(df > 25) & (df < 35)`). Although the speed difference may appear negligible on smaller datasets, the performance gains become substantial and measurable when processing datasets involving millions of rows.

Clarity and Readability of Code: Perhaps the most significant non-technical advantage of integrating `.between()` into your coding habits is the immediate and profound enhancement in code readability. The method clearly and concisely communicates the exact intent of selecting values within a specific range, including the precise boundary conditions defined by the `inclusive` parameter. This structural clarity dramatically simplifies the debugging process, minimizes the opportunity for logical errors inherent in complex conditional statements, and significantly improves overall code maintainability for future collaboration.

Conclusion

The `.between()` method within [Pandas](#) offers a powerful, efficient, and highly readable solution for handling the fundamental data manipulation task of filtering DataFrame rows based on range constraints. Its inherent flexibility, largely enabled by the granular control provided through the `inclusive` parameter, permits analysts to rapidly define precise open, closed, or semi-open intervals, thereby accommodating diverse and specific analytical demands with minimal, expressive code. Furthermore, the capacity to straightforwardly invert the selection using the [tilde](#) operator (`~`) provides an elegant and concise mechanism for anomaly detection and the isolation of data points falling outside expected operational thresholds.

By effectively integrating the `.between()` method into your analytical toolkit, you can significantly streamline and clarify your data filtering workflows, making your [data science](#) and analysis tasks substantially more accurate and efficient. Always maintain the best practice of confirming the data types of your columns and critically understanding the implications for [missing values \(NaN\)](#) to ensure that your filtering operations consistently yield reliable and trustworthy results.

Additional Resources

To further enhance your Pandas proficiency and explore more sophisticated data manipulation and filtering techniques, the following resources are recommended:

[Pandas User Guide on Indexing and Selecting Data](#): A comprehensive and authoritative guide to various data selection methods in Pandas.

[Pandas `.isin\(\)` Method Documentation](#): Learn how to select DataFrame rows where a column's value is present in a list of specified values.

[Comprehensive Guide to Filtering DataFrames in Pandas](#): An external tutorial offering a broader perspective on various DataFrame filtering strategies and techniques.