

How to Clean Your Data with Pandas dropna() and Thresh

Authored by
stats writer

November 20, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Clean Your Data with Pandas dropna() and Thresh*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=98536>

The process of Pandas data cleaning often requires robust methods for handling missing values. When dealing with real-world datasets, the presence of these gaps--typically represented as Not a Number (NaN)--can severely impact statistical analyses and machine learning model performance. One of the most critical tools for addressing this challenge is the `dropna()` function, which provides flexible ways to eliminate rows or columns based on their null status.

Specifically, the `dropna()` function, when combined with the **`thresh`** parameter, allows data practitioners to impose a minimum survival criterion. This powerful feature lets you specify the minimum number of non-null observations required for a given row or column to be retained in the resulting `DataFrame`. This approach moves beyond simple removal of any row containing *any* missing value, offering a nuanced and controlled mechanism for data imputation and preparation.

The Power of `dropna()`: Basic Usage and Control

At its core, the `dropna()` function is designed to systematically drop rows or columns from a `DataFrame` that contain missing values. While the default behavior (`how='any'`) removes any row with at least one null entry, this can often lead to excessive data loss, especially in sparse datasets where missingness is common but isolated.

To mitigate aggressive data deletion, Pandas provides the **`thresh`** argument. This argument fundamentally changes the function's logic from a binary check ("Is there a NaN?") to a quantitative assessment ("Are there enough non-NaN values?"). By setting **`thresh`**, you are defining the minimal threshold of acceptable data fidelity for an observation (row) or a feature (column) to be preserved.

Understanding the interplay between **`thresh`** and the **`axis`** parameter is essential. When `axis=0` (the default, operating on rows), **`thresh`** specifies the minimum count of non-null cells within that row. Conversely, setting `axis=1` instructs the function to evaluate columns, and **`thresh`** then dictates the minimum number of non-null observations required within that column to keep it in the dataset. This flexibility ensures that you can target data quality requirements precisely where they are needed.

Mastering the `thresh` Parameter for Precision

The **`thresh`** argument provides granular control over how missing data is handled, moving beyond simple deletion strategies. It is particularly valuable when dealing with large datasets where some level of missingness is expected but complete removal would cause the loss of valuable, non-null information contained elsewhere in the record.

When defining the threshold value, it is crucial to consider the context of your data and the potential impact of data density on your downstream analysis. The value passed to **`thresh`** must be

an integer, representing the absolute count of non-null entries. For instance, if a DataFrame has 10 columns, setting `thresh=8` (on rows) means that only rows possessing 8 or more valid (non-NaN) data points will be retained.

The following four methods illustrate the practical application of the **thresh** argument, showcasing how to use absolute counts versus calculated percentages for both row-wise and column-wise filtering. These methods represent the most common and powerful ways to leverage this parameter for sophisticated missing data management.

Method 1: Only Keep Rows with Minimum Number of non-NaN Values

```
#only keep rows with at least 2 non-NaN values  
df.dropna(thresh=2)
```

Method 2: Only Keep Rows with Minimum % of non-NaN Values

```
#only keep rows with at least 70% non-NaN values  
df.dropna(thresh=0.7*len(df.columns))
```

Method 3: Only Keep Columns with Minimum Number of non-NaN Values

```
#only keep columns with at least 6 non-NaN values  
df.dropna(thresh=6, axis=1)
```

Method 4: Only Keep Columns with Minimum % of non-NaN Values

```
#only keep columns with at least 70% non-NaN values  
df.dropna(thresh=0.7*len(df), axis=1)
```

To demonstrate these methods effectively, we will utilize a sample dataset representing sports statistics. This dataset includes various columns (team, points, assists, rebounds) and contains deliberate missing values, allowing us to clearly observe the filtering effects of the **thresh** parameter across different configurations.

```
import pandas as pd  
import numpy as np
```

```
#create DataFrame  
df = pd.DataFrame({'team': ,  
'points': ,
```

```
'assists': ,
'rebounds': })

#view DataFrame
print(df)

team points assists rebounds
0 A 18.0 5.0 11.0
1 B NaN NaN NaN
2 C 19.0 NaN 10.0
3 D 14.0 9.0 6.0
4 E 14.0 NaN 6.0
5 F 11.0 9.0 5.0
6 G 20.0 9.0 9.0
7 H NaN 4.0 NaN
```

Example 1: Filtering Rows by Absolute Non-Null Count

When we are confident about the minimum amount of information required for a row to be considered useful, we use the absolute count method. In this example, our `DataFrame` has four columns. We decide that any player record must have at least two valid data points (non-`NaN` values) to justify keeping the row. This is implemented by setting `thresh=2` without specifying the `axis` parameter, thereby defaulting to row operations.

This approach ensures that rows with extreme sparsity, such as those that are entirely or nearly entirely null, are efficiently discarded. It is a fundamental strategy for initial data quality control, ensuring that subsequent analyses are performed only on records with a minimum level of integrity.

We apply the following syntax to retain only those rows possessing a minimum of two non-`NaN` values:

```
#only keep rows with at least 2 non-NaN values
df.dropna(thresh=2)
```

```
team points assists rebounds
0 A 18.0 5.0 11.0
2 C 19.0 NaN 10.0
3 D 14.0 9.0 6.0
4 E 14.0 NaN 6.0
5 F 11.0 9.0 5.0
```

```
6 G 20.0 9.0 9.0
7 H NaN 4.0 NaN
```

Upon execution, we observe that the row corresponding to index position **1** (Team B) has been successfully dropped from the `DataFrame`. This outcome occurs because Team B's record contained only one non-null entry (the 'team' name itself), failing to meet the specified `thresh=2` requirement. All other rows contained two or more valid values and were thus retained.

Example 2: Filtering Rows by Percentage Threshold

While using absolute counts is useful, dataframes often change in size or feature count. A more dynamic and often preferred method for filtering rows is using a percentage threshold, which calculates the required non-null count based on the current number of columns. This provides a robust solution that adapts even if the dataset schema is modified later.

In this scenario, we decide that a row must have at least 70% of its values present to be kept. Since our sample `DataFrame` has four columns, 70% translates to $0.7 * 4 = 2.8$. Since the `thresh` argument must be an integer, Pandas effectively rounds this up (or takes the ceiling of the calculation), meaning we require a minimum of 3 non-null values per row.

We calculate the required threshold dynamically using the length of the columns and apply the `thresh` parameter as follows:

```
#only keep rows with at least 70% non-NaN values
df.dropna(thresh=0.7*len(df.columns))
```

```
team points assists rebounds
0 A 18.0 5.0 11.0
2 C 19.0 NaN 10.0
3 D 14.0 9.0 6.0
4 E 14.0 NaN 6.0
5 F 11.0 9.0 5.0
6 G 20.0 9.0 9.0
```

Comparing this output to Example 1, we note that both index position **1** (Team B, 1 non-null) and index position **7** (Team H, 2 non-null) have been dropped. Team H was retained in Example 1, but failed here because it only had 2 non-null values, falling short of the required 3 points (70% of 4 columns). This highlights how setting a stricter percentage-based threshold leads to higher data density.

Example 3: Filtering Columns by Absolute Non-Null Count

While row filtering addresses records with insufficient information, column filtering targets features (variables) that are too sparse to be predictive or reliable. When performing column-wise deletion, we must explicitly set the `axis` parameter to `1`. The value provided to `thresh` then represents the minimum required number of non-null rows needed for a column to survive the operation.

For our running example, we determine that any column must contain at least six valid observations out of eight total rows to be considered valuable for analysis. We specify `thresh=6` and `axis=1`.

This method is paramount in feature engineering and selection, preventing the inclusion of columns that are dominated by NaNs, which often introduce noise or require complex imputation techniques if retained.

#only keep columns with at least 6 non-NaN values

`df.dropna(thresh=6, axis=1)`

team points rebounds

0 A 18.0 11.0

1 B NaN NaN

2 C 19.0 10.0

3 D 14.0 6.0

4 E 14.0 6.0

5 F 11.0 5.0

6 G 20.0 9.0

7 H NaN NaN

In the resulting output, the 'assists' column has been removed. Reviewing the original data confirms that 'assists' only contained five non-null entries (5, 9, 9, 9, 4), falling below our required threshold of six. Conversely, the 'team', 'points', and 'rebounds' columns all met or exceeded this count and were preserved.

Example 4: Filtering Columns by Percentage Threshold

Similar to row filtering, column filtering can be defined using a percentage of the total dataset length, ensuring scalability and consistency regardless of the total number of records. If the size of the dataset changes (e.g., new data is appended), the threshold automatically adjusts relative to the current number of rows.

We will again set a requirement that columns must contain at least 70% non-NaN values. Given

that the sample DataFrame has eight rows, 70% translates to $0.7 * 8 = 5.6$. This means any column must contain a minimum of 6 non-null values to be retained.

Using a percentage threshold dynamically is often considered best practice when the integrity of features is paramount, ensuring that columns are only used if they represent a majority of the observations.

#only keep columns with at least 70% non-NaN values

`df.dropna(thresh=0.7*len(df), axis=1)`

```
team points rebounds
```

```
0 A 18.0 11.0
```

```
1 B NaN NaN
```

```
2 C 19.0 10.0
```

```
3 D 14.0 6.0
```

```
4 E 14.0 6.0
```

```
5 F 11.0 5.0
```

```
6 G 20.0 9.0
```

```
7 H NaN NaN
```

As seen in the output, the 'assists' column is again dropped. Since this column only contained five valid entries, it failed to meet the calculated requirement of six non-NaN values (70% of 8 rows). The result is identical to Example 3, but demonstrates the flexibility of using percentage calculations for setting the **thresh** value dynamically.

Conclusion: Best Practices for Data Cleaning

The use of the **thresh** parameter within the `dropna()` function is an indispensable technique for robust data preparation in Pandas. It allows data scientists to move beyond simplistic strategies for missing data handling, providing a precise mechanism for quality control based on data density.

Choosing the correct threshold--whether absolute or percentage-based--should always be guided by domain knowledge and the requirements of the subsequent analytical task. Aggressive filtering can lead to data loss, while overly lenient filtering may introduce noise. The methods demonstrated here provide the fundamental tools for navigating this balance effectively.

For detailed information regarding all parameters and advanced configurations of this powerful method, consult the official [pandas](#) documentation. Implementing these controlled data cleaning techniques ensures that your datasets maintain high fidelity for accurate modeling and reporting.

Note: You can find the complete documentation for the pandas **dropna()** function [here](#).