

# Learning to Identify Missing Data: A Guide to Using “Is Not Null” in Pandas

Authored by  
**Mohammed loot**

May 28, 2026

## RECOMMENDED CITATION

Mohammed loot (2026). *Learning to Identify Missing Data: A Guide to Using “Is Not Null” in Pandas*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=3665>

In the complex process of [data analysis](#) and manipulation, particularly when leveraging the power of [Pandas](#), mastering the handling of missing data is absolutely critical. These gaps, frequently represented as the floating-point value **NaN** (Not a Number) or Python's built-in constant [None](#), can severely compromise the integrity and reliability of any statistical or analytical output. To efficiently combat this challenge, Pandas offers an array of robust methods, with the [notnull\(\)](#) function standing out as a fundamental utility for identifying valid observations.

The [notnull\(\)](#) function provides a straightforward and highly effective means of testing for the presence of sound, non-missing data across an entire [Pandas DataFrame](#) or a single [Series](#). When executed, it meticulously evaluates every single element within the structure and generates a corresponding [boolean](#) result. A return of **True** confirms that the element holds a discernible, non-null value, while **False** indicates the element is null (i.e., it is **NaN** or [None](#)). This logical structure, which is the direct inverse of the more commonly cited [isnull\(\)](#) function, is especially intuitive when the primary goal is to isolate and focus solely on the complete portions of the dataset.

A comprehensive understanding of how to effectively utilize [notnull\(\)](#) is paramount for effective [data cleaning](#) and preparing raw datasets for subsequent sophisticated processing. Whether your objective involves systematically filtering out rows that are incomplete, rapidly identifying which columns suffer from significant missing entries, or simply generating a precise count of all valid observations, this function provides the essential groundwork. The following sections are dedicated to exploring several highly practical and common applications of [notnull\(\)](#), utilizing clear, step-by-step examples that vividly illustrate its versatility across diverse data manipulation challenges.

## Establishing the Foundation: Creating Our Example DataFrame

To properly illustrate the full range of capabilities offered by the [notnull\(\)](#) function, our initial step involves the construction of a representative sample [Pandas DataFrame](#). This synthetic dataset is engineered to realistically mimic real-world data environments by including a deliberate mixture of numerical and categorical data types, alongside strategically positioned null values. This setup ensures we have a tangible example where missing information is a known factor. We rely on the core [Pandas](#) library for the DataFrame construction itself, and integrate the [NumPy](#) library specifically for the efficient representation of **NaN** values.

Our resulting DataFrame is designed to hold a small, illustrative collection of sports team statistics, encompassing common metrics such as points scored, assists recorded, and rebounds collected. By intentionally embedding these missing data points, we create a clear and observable scenario where we can demonstrate precisely how [notnull\(\)](#) functions to identify and enable effective management of these absent data entries. This realistic context is vital, as it allows us to see the practical impact and utility of the function in action for every subsequent demonstration.

Presented below is the complete Python code utilized to generate our example DataFrame. Following the creation script, you will find the printed output of the DataFrame, which will serve as the consistent baseline dataset for all the filtering and aggregation examples that follow:

```
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 22.0 NaN 8.0
2 C 19.0 7.0 10.0
3 D 14.0 9.0 6.0
4 E 14.0 12.0 6.0
5 F 11.0 9.0 5.0
6 G 20.0 9.0 NaN
7 H NaN NaN 12.0
```

## Example 1: Isolating Complete Observations Across All Columns

A frequent and often mandatory requirement in rigorous [data cleaning](#) processes is the ability to swiftly identify and retain only those rows that are entirely free of missing values across every single column. This highly selective filtering methodology guarantees that every observation included in the subsequent analysis is complete and fully intact, thereby preemptively eliminating potential distortions or inaccuracies that are typically introduced by processing incomplete data points. The powerful [notnull\(\)](#) function, when intelligently combined with other core [Pandas](#) aggregation methods, makes this crucial filtering operation both straightforward to implement and highly efficient in execution.

To successfully achieve this comprehensive row completeness filter, we initiate the process by applying [notnull\(\)](#) to our entire [Pandas DataFrame](#). This action immediately yields a new [boolean](#) DataFrame that possesses the exact same dimensions as the original structure. In this resultant

DataFrame, each cell is marked **True** if the corresponding value in the original cell was non-null, and **False** if a null value was present. Crucially, we then chain the **.all(1)** method onto this boolean result. The dedicated **.all()** method is designed to check if every element along a specific axis evaluates to **True**. By setting **axis=1** (which specifies row-wise operation), the method returns a specialized [Series](#) of **True** for rows where all values are non-null, and **False** for any row that contains even a single null value.

This derived [boolean Series](#) is then seamlessly employed to index our starting [Pandas DataFrame](#), thereby enabling a powerful filter that displays only those rows that have zero null values across any of their columns. This specific methodology proves exceptionally valuable in contexts where absolute data completeness is a strict prerequisite for any subsequent analytical tasks or modeling efforts.

The following precise code snippet demonstrates the process of filtering the [DataFrame](#) to retain only those rows where every single column contains a valid, non-null data entry:

#### **#filter for rows with no null values in any column**

**df**

```
team points assists rebounds
```

```
0 A 18.0 5.0 11.0
```

```
2 C 19.0 7.0 10.0
```

```
3 D 14.0 9.0 6.0
```

```
4 E 14.0 12.0 6.0
```

```
5 F 11.0 9.0 5.0
```

As clearly evident from the resulting output, rows 1, 6, and 7 from the initial [DataFrame](#) have been systematically excluded because each contained a minimum of one [NaN](#) value. Consequently, the resulting [DataFrame](#) now exclusively features rows where all critical entries--'team', 'points', 'assists', and 'rebounds'--are entirely complete.

### **Example 2: Granular Filtering: Ensuring Completeness in a Specific Column**

In numerous [data analysis](#) situations, the required level of data completeness might be narrowly focused on one or two particular columns, rather than necessitating completeness across the entire row. For instance, if a specific variable is central to a critical calculation, it becomes paramount to ensure that all observations used in that calculation possess a valid and non-missing value for that single metric. The [notnull\(\)](#) function excels in this precise targeting, offering users highly granular and flexible control over their filtering criteria.

To effectively filter rows based strictly on the non-null status of a single column, the procedure

begins by first explicitly selecting the column of interest (often done by passing a list containing the column name, e.g., `df`). Applying `notnull()` to this focused selection immediately generates a dedicated [boolean Series](#), which accurately flags which elements in that specified column are non-null. Following this, we chain the `.all(1)` method, which, while simple for a single column, maintains consistent syntax and confirms that the required value is present in that column for each row. The output of this operation is the final [boolean Series](#) used to index and subset the main DataFrame.

This highly focused method enables practitioners to concentrate their limited [data cleaning](#) resources precisely on the most critical features within their dataset. Taking the example of our statistics, if the 'assists' column is deemed indispensable for a particular analysis, we can selectively filter out all rows where this specific value is absent, entirely disregarding the completeness status of other, less critical columns like 'points' or 'rebounds'. This adaptability is an invaluable asset for preserving data integrity exactly where it holds the most significance for the predefined analytical goals.

The following implementation demonstrates how to precisely filter the [DataFrame](#) to include only those rows that contain a valid value specifically within the 'assists' column:

```
#filter for rows with no null values in the 'assists' column  
df[notnull().all(1)]
```

```
team points assists rebounds  
0 A 18.0 5.0 11.0  
2 C 19.0 7.0 10.0  
3 D 14.0 9.0 6.0  
4 E 14.0 12.0 6.0  
5 F 11.0 9.0 5.0  
6 G 20.0 9.0 NaN
```

It is important to observe that in this newly filtered [DataFrame](#), every row now consistently features a valid, non-null value in the 'assists' column. Crucially, rows such as 'G' are successfully retained, even though their 'rebounds' column contains a [NaN](#) value, because the filtering operation was focused entirely and exclusively on the 'assists' column. This example clearly showcases the precision and efficacy of targeted filtering using `notnull()`.

### Example 3: Quantitative Assessment: Counting Valid Entries Per Column

Moving beyond simple filtering, gaining a clear quantitative understanding of the distribution of missing data across the [DataFrame](#)'s columns is essential for accurately assessing the overall data quality and effectively pinpointing areas that require immediate [data cleaning](#) intervention.

Counting the precise number of non-null values within each column delivers a rapid, high-level overview of how complete each feature is, which in turn informs strategic decisions regarding imputation techniques, subsequent feature selection processes, or even the necessity for improved data collection protocols.

The `notnull()` function, when powerfully combined with the aggregation capabilities of the `.sum()` method, provides an exceptionally clean and direct pathway to obtaining these vital counts. When `notnull()` is initially applied to any `DataFrame`, its output is a `boolean` `DataFrame` where the value `True` signifies a non-null entry and `False` denotes a null entry. Within Python's standard arithmetic operations, `True` is inherently interpreted as the integer 1, and `False` is interpreted as 0.

Therefore, invoking `.sum()` on this resulting `boolean` `DataFrame` effectively tallies the total number of `True` values (which correspond exactly to the non-null entries) for every column. By default, `.sum()` operates in a column-wise manner (using `axis=0`), generating a `Series` where each individual entry represents the exact count of non-null values for its corresponding column. This aggregated perspective is invaluable for quickly isolating which specific columns might be suffering from significant and unacceptable data gaps.

The following code demonstrates the methodology used to accurately count the number of non-null values currently present within each column of our example `DataFrame`:

```
#count number of non-null values in each column  
df.notnull().sum()
```

```
team 8  
points 7  
assists 6  
rebounds 7  
dtype: int64
```

From this detailed output, we gain immediate and actionable insights into the completeness status of each feature column:

The **team** column has 8 non-null values, confirming it is completely filled with data.

The **points** column has 7 non-null values, indicating that exactly one entry is currently missing.

The **assists** column has 6 non-null values, signifying that two entries are absent.

The **rebounds** column has 7 non-null values, showing that one entry is missing.

This itemized breakdown by column is fundamentally invaluable for conducting preliminary data quality checks and initiating targeted `data cleaning` efforts.

## Example 4: Calculating Total Data Density: Non-Null Count Across the Entire DataFrame

While obtaining the non-null counts on a column-by-column basis provides necessary detail, there are specific contexts--such as high-level data quality reporting or pre-processing comparisons--where a single, aggregated number representing the total count of all valid data points across the entirety of the [DataFrame](#) is required. This single metric offers a powerful summary of the dataset's density. This overall count is extremely useful for rapid data quality assurance checks, for systematically comparing the total volume of valid data before and after significant [data cleaning](#) operations, or for concise reporting to stakeholders. The `notnull()` function, when paired with a sequential double summation, presents an extremely elegant and efficient solution for calculating this grand total.

Building directly upon the methodology established in the previous example, we initiate the process by applying `notnull()` to our [DataFrame](#), yielding the familiar [boolean](#) DataFrame. The first application of `.sum()` effectively aggregates all the **True** values (the column-wise non-null counts) for each individual column, which results in the production of a [Series](#) object. This Series now contains the total valid entries for every feature.

To finally arrive at the grand total count spanning the entire [DataFrame](#), we simply apply the `.sum()` method a second time, this time to the resulting [Series](#) itself. This second summation process adds up all the previously calculated column-wise non-null counts, culminating in a single integer scalar value. This scalar concisely represents the overall count of all valid entries contained within the entire dataset, offering a crucial metric for overall data completeness.

The following code snippet clearly illustrates the technique used to calculate the absolute total number of non-null values aggregated across all columns and all rows within the example [DataFrame](#):

```
#count number of non-null values in entire DataFrame  
df.notnull().sum().sum()
```

28

Based on this definitive output, we can conclusively state that there are exactly **28** non-null values distributed across the entirety of the [DataFrame](#). This single, powerful number provides an instantaneous summary of the data density, enabling rapid and informed assessment of the dataset's overall quality and completeness at a single glance.

## Conclusion and Recommended Next Steps

The [Pandas `notnull\(\)`](#) function is correctly regarded as an indispensable and foundational tool for any individual engaged in data manipulation using Python. As comprehensively demonstrated through the practical examples provided, it offers a versatile, efficient, and highly precise mechanism to identify, filter, and accurately quantify all non-missing values contained within your [DataFrames](#). From the targeted filtering of complete rows to the quantitative assessment of column-wise or overall data completeness, [notnull\(\)](#) critically empowers users to make well-informed decisions regarding essential [data cleaning](#), preprocessing methodologies, and subsequent [data analysis](#) activities.

Achieving proficiency with this function represents a fundamental and necessary milestone toward the construction of robust, efficient, and ultimately reliable data processing pipelines. Its intuitive [boolean](#) output, when seamlessly integrated with [Pandas](#)' formidable indexing and aggregation capabilities, solidifies its position as a cornerstone technique in the effective management of missing data. By consciously integrating these techniques into your standard data workflow, you can significantly elevate both the quality and integrity of your resulting datasets, ensuring more trustworthy analytical outcomes.

For those aspiring to further deepen their expertise in the comprehensive management of missing data within [Pandas](#), it is highly recommended to explore several related and complementary functions. These include [.isnull\(\)](#) (which operates as the direct logical inverse of [notnull\(\)](#)), [.dropna\(\)](#) for the systematic removal of rows or columns containing missing data, and [.fillna\(\)](#) for implementing strategies to impute or fill in those missing values.

The following tutorials offer further explanations on how to execute other common and necessary filtering operations within the [Pandas](#) environment: