

# Fill NA Values for Multiple Columns in Pandas

Authored by  
**Mohammed loot**

November 6, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Fill NA Values for Multiple Columns in Pandas*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=11286>

## The Crucial Importance of Managing Missing Data

In the expansive fields of **data analysis** and **machine learning**, encountering [missing values](#) is not merely a possibility--it is an inevitable challenge. These data voids, frequently denoted as Not a Number (NaN) or null markers, possess the capacity to severely compromise the integrity of statistical assessments, derail the efficacy of model training, and ultimately lead to inaccurate or misleading conclusions. Consequently, the proficient and thoughtful management of these omissions constitutes one of the most fundamental and critical steps within the entire **data preparation pipeline**.

The [Pandas](#) library, recognized globally as a foundational pillar of the Python data science ecosystem, furnishes a robust and highly adaptable solution for tackling this specific difficulty: the dedicated [fillna\(\)](#) function. This versatile function empowers data practitioners to systematically replace missing data points using a comprehensive range of strategies, spanning from the simplicity of substituting with a basic constant to the sophistication of advanced statistical **interpolation techniques**.

This detailed guide is engineered to explore and illuminate several highly effective techniques for strategically employing the `fillna()` method, particularly when the objective involves simultaneously addressing missing data across multiple columns within a single [Pandas DataFrame](#). By mastering these specialized methods, analysts can ensure unparalleled **data integrity** and successfully prepare their datasets for the demands of rigorous and reliable analysis.

## Establishing the Environment and Constructing the Sample Dataset

Before proceeding with the practical application of **imputation strategies**, it is essential to properly configure our working environment. This involves importing the necessary core libraries: [Pandas](#), for data manipulation, and [NumPy](#), which is crucial as it provides the standard representation for [missing values](#) (`np.nan`) used throughout numerical arrays and columns in Python. Understanding the role of `np.nan` is key to accurately identifying and targeting null entries.

To provide a clear context for our subsequent examples, we will utilize a carefully constructed sample [DataFrame](#). This dataset is designed to realistically simulate common sports statistics, incorporating a mix of data types--specifically, categorical strings and continuous numerical values--alongside several strategically embedded null entries across various columns. This setup allows us to accurately demonstrate the necessary flexibility required when cleaning and preparing truly complex, real-world data structures.

The following code block executes the creation of our initial sample [DataFrame](#) and then prints its structure to the console. The output clearly highlights the specific locations of the null entries that must be targeted and addressed using the different methods of the `fillna()` function:

```
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 25.0 5.0 11
1 NaN NaN 7.0 8
2 B 15.0 7.0 10
3 B NaN 9.0 6
4 B 19.0 12.0 6
5 C 23.0 9.0 5
6 C 25.0 NaN 9
7 C 29.0 4.0 12
```

## Method 1: Universal Imputation Using a Single Constant

The most straightforward implementation of the [Pandas fillna\(\)](#) function is known as **global imputation**. In this technique, a singular, unchanging constant value is applied uniformly to substitute all [missing values](#) throughout the entirety of the [DataFrame](#). While this method offers immediate results and simplicity, it requires significant caution. Imputing the same value across columns that represent fundamentally different entities or metrics (for example, applying '0' to both a categorical team name and a numerical points tally) carries a high risk of introducing substantial statistical bias and data type inconsistencies.

For the purpose of this demonstration, we elect to replace every missing entry, irrespective of its original column, with the numerical value **zero (0)**. This specific approach is sometimes deemed acceptable when working strictly with count data or quantitative variables where a null entry can logically be interpreted as a genuine absence of the measured quantity. We leverage the `inplace=True` parameter, which efficiently modifies the [DataFrame](#) directly, thereby eliminating the need for an explicit assignment step.

The following code demonstrates this universal replacement strategy. It is essential to observe how

the NaN markers in the 'team' (originally a string/object type), 'points' (float), and 'assists' (float) columns are all simultaneously replaced by '0'. This action often results in the forced **type coercion** of columns, a common side effect of global imputation that requires careful post-processing:

```
#replace all missing values with zero
```

```
df.fillna(value=0, inplace=True)
```

```
#view DataFrame
```

```
print(df)
```

```
team points assists rebounds
```

```
0 A 25.0 5.0 11
```

```
1 0 0.0 7.0 8
```

```
2 B 15.0 7.0 10
```

```
3 B 0.0 9.0 6
```

```
4 B 19.0 12.0 6
```

```
5 C 23.0 9.0 5
```

```
6 C 25.0 0.0 9
```

```
7 C 29.0 4.0 12
```

While extremely effective for rapid data cleanup, analysts must consistently and critically evaluate the **semantic implications** of utilizing zero imputation. If a column represents a characteristic where zero is an inherently improbable, impossible, or statistically misleading value (such as age, income, or a complex measurement), then a blanket constant replacement is likely an inappropriate and detrimental analytical choice.

## Method 2: Targeted Imputation for Specific Columns

A significantly more precise and generally preferred methodology involves applying imputation selectively by explicitly targeting only a subset of columns. This strategy ensures that the designed replacement logic is executed only in the locations where it is contextually and statistically sound, thereby preserving the integrity and original data types of all untouched columns. This targeted technique is particularly invaluable when the goal is to clean up a limited number of numerical or categorical features without affecting the rest of the dataset.

To implement this method effectively, we utilize standard **list notation** to select the necessary columns, which generates a temporary view of that specific subset of the [DataFrame](#). We then apply the [fillna\(\)](#) method directly to this selection. Crucially, the imputed results must then be explicitly assigned back to the original subset of the DataFrame to finalize the modification. This

process ensures only the specified columns are updated.

In the following instance, we choose to impute only the 'points' and 'assists' columns, replacing their [NaN](#) entries with the value zero. Note that before executing this step, we assume the use of a fresh copy of the original DataFrame (or that the previous global modifications have been undone). As intended, this targeted approach leaves the missing entries within the 'team' column entirely untouched, demonstrating precise control over the cleaning process:

```
#replace missing values in points and assists columns with zero
```

```
df] = df].fillna(value=0)
```

```
#view DataFrame
```

```
print(df)
```

```
team points assists rebounds
0 A 25.0 5.0 11
1 NaN 0.0 7.0 8
2 B 15.0 7.0 10
3 B 0.0 9.0 6
4 B 19.0 12.0 6
5 C 23.0 9.0 5
6 C 25.0 0.0 9
7 C 29.0 4.0 12
```

The inherent flexibility of this technique is highly advantageous. For example, instead of using a constant zero, one could replace missing values in 'points' with the calculated mean of the 'points' column, while simultaneously replacing missing values in 'assists' with the median of the 'assists' column. Such customized operations can be elegantly achieved by applying calculated values within the `value` parameter or through **method chaining**.

### Method 3: Granular Control Using Dictionary Mapping

When working with complex datasets where multiple columns necessitate entirely unique imputation values or distinct statistical strategies, the utilization of a Python **dictionary** within the [fillna\(\)](#) function provides the ultimate level of granular control. This advanced method permits the direct mapping of specific column names (acting as dictionary keys) to their corresponding, tailor-made replacement values or functions (acting as dictionary values).

The syntax for implementing this powerful method is highly intuitive: `df.fillna({'column_name_1': value_1, 'column_name_2': value_2, ...})`. This structure is essential for successfully handling heterogeneous data types--allowing us to treat categorical

columns separately from numerical ones, thereby ensuring that data types are either properly maintained or intentionally modified to suit the analytical requirement.

In our conclusive example of the primary methods, we illustrate how to replace [NaN](#) entries across three different columns, each utilizing a distinct replacement value tailored to its data type or intended meaning:

The categorical 'team' column is strategically filled with the string '**Unknown**'.

The numerical 'points' column is filled using the integer **0**.

The numerical 'assists' column is deliberately filled with the string '**zero**'. (A critical observation here: this intentional use of a string forces the entire 'assists' column to undergo **type conversion** to an object/string data type.)

**#replace missing values in three columns with three different values**

```
df.fillna({'team':'Unknown', 'points': 0, 'assists': 'zero'}, inplace=True)
```

```
#view DataFrame
```

```
print(df)
```

```
team points assists rebounds
```

```
0 A 25.0 5 11
```

```
1 Unknown 0.0 7 8
```

```
2 B 15.0 7 10
```

```
3 B 0.0 9 6
```

```
4 B 19.0 12 6
```

```
5 C 23.0 9 5
```

```
6 C 25.0 zero 9
```

```
7 C 29.0 4 12
```

As clearly demonstrated by the output, the missing entries are replaced precisely according to the defined dictionary mapping. The 'team' entry is correctly updated to 'Unknown'. Furthermore, the 'assists' column now contains a mixture of numerical values and the string literal 'zero', which definitively confirms that [Pandas](#) has coerced its underlying **data type** to 'object' to successfully accommodate the heterogeneous values. This outcome serves as a potent reminder of the importance of thoroughly understanding the data type implications inherent in any complex imputation strategy.

## Strategic Considerations for Robust Imputation

While the mechanical process of replacing [missing values](#) with a constant (such as 0 or 'Unknown') is straightforward, it is rarely considered the optimal analytical decision. The selection of the correct

imputation strategy is critically important for minimizing statistical bias, ensuring the fidelity of the dataset, and maintaining the integrity of all subsequent analyses. The three methods detailed previously provide the structural mechanisms for applying the chosen value; however, the value itself must be determined through rigorous **statistical justification**.

Data professionals routinely employ several sophisticated strategies and best practices when utilizing [fillna\(\)](#) to achieve superior results. These advanced techniques go beyond simple constant replacement and address the underlying statistical properties of the data:

**Statistical Imputation (Mean, Median, Mode):** For columns containing **numerical data**, replacing missing points with the column's mean or median is standard practice. The [DataFrame](#) methods `.mean()`, `.median()`, and `.mode()` can be highly effective when integrated directly into the dictionary method (Method 3). The median is often the preferred choice if the data distribution exhibits heavy skewness, as it is significantly less sensitive to the influence of **outliers** than the mean.

**Handling Categorical Data:** For columns that contain categorical variables (like our 'team' column), the common strategies involve replacing NaN entries with the mode (the most frequently occurring category) or assigning a specific, descriptive label such as 'Missing' or 'Unknown'. Using such a placeholder label is beneficial because it explicitly preserves the information that the data was missing, treating it as a distinct category that can be valuable for downstream **predictive modeling**.

**Temporal and Sequential Imputation:** When the dataset involves **time-series data** or other sequentially ordered observations, basic statistical averages are often insufficient. [Pandas](#) offers specialized methods crucial for maintaining temporal relationships, including **Forward Fill** (`ffill`, which propagates the last observed non-null value forward) and **Backward Fill** (`bfill`, which uses the next observed non-null value). These are essential for preserving the flow of sequential data.

**Group-Based Conditional Imputation:** One of the most analytically powerful strategies involves calculating the replacement value based on specific subgroups within the data rather than the overall population. For instance, instead of filling a missing 'points' value with the general mean score, you would replace it with the mean points scored *\*specifically by that team\**. Implementing this requires the sophisticated use of the `groupby()` method in conjunction with `transform()` before the final application of [fillna\(\)](#).

Finally, regardless of the method chosen, maintaining meticulous **documentation** is paramount. Always record which imputation strategy was deployed for each column. Furthermore, it is considered best practice to conduct **sensitivity testing** to evaluate how robust your final statistical model or analysis is to the chosen imputation method, ensuring that the results are reliable and not artifacts of the cleaning process.

## Conclusion and Future Directions

The `fillna()` function stands as an indispensable tool within the data preparation workflow in Python, providing essential mechanisms for reliably managing the ubiquitous challenge of missing data. Throughout this guide, we have systematically covered the three fundamental methods for deploying this function across multiple columns in a [Pandas DataFrame](#):

**Global Replacement:** Applying a single, uniform constant across the entire dataset.

**Targeted Replacement:** Using list notation to select and impute specific columns.

**Advanced Dictionary Mapping:** Utilizing a dictionary to achieve highly granular control by assigning unique values or strategies to individual columns.

By achieving mastery over these implementation mechanisms, you gain the necessary foundation to select and execute statistically sound imputation strategies that are tailored to the nuances of your dataset. Effective data cleaning is undeniably the first and most critical prerequisite for constructing reliable predictive models and accurately deriving insights from complex, real-world data structures. We strongly encourage further exploration into highly advanced techniques, such as statistical modeling for missing data, iterative imputation methods, or the use of specialized machine learning models designed explicitly to predict and fill [NaN](#) values, pushing the boundaries of data integrity.