

Learning to Apply Functions to Multiple Columns in Pandas DataFrames

Authored by
Mohammed loot

November 13, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning to Apply Functions to Multiple Columns in Pandas DataFrames*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=23988>

When conducting sophisticated data analysis on substantial datasets using the [Pandas](#) library in Python, data scientists frequently encounter scenarios where standard, built-in functions are inadequate for complex data transformation needs. Often, the requirement is to define a custom, nuanced logic that operates on the values across multiple columns simultaneously within a single observation, or [DataFrame](#) row. Although Pandas is highly optimized for [vectorization](#)--operations that apply uniformly across entire columns--applying complex, row-wise calculations that depend on several input variables demands a more flexible and iterative approach. This is precisely where the powerful and essential combination of the **apply()** method and a **lambda** expression proves indispensable.

Mastering this technique unlocks highly flexible data manipulation capabilities, allowing users to calculate complex derived metrics, implement conditional transformations based on multiple criteria, or execute intricate mathematical models where inputs are sourced from various fields within the same record. For efficient and idiomatic data processing in Python, understanding how to correctly structure this application, particularly the critical role played by the **axis** parameter, is absolutely fundamental.

Leveraging `df.apply()` with Anonymous `lambda` Functions

The standardized and most accessible method for applying custom logic across multiple columns of a [DataFrame](#) involves utilizing the [apply\(\)](#) method in conjunction with an anonymous [lambda function](#). The primary purpose of the **apply()** method is to execute a specified function along either the index (rows) or columns axis of a DataFrame, establishing it as a highly versatile component of the data transformation toolkit. When the goal is to access and process the data from several columns corresponding to a single record, the [lambda function](#) provides the necessary mechanism to define the operation inline, directly bridging the DataFrame structure to the custom logic.

The [lambda function](#) serves as a lightweight wrapper, taking the entire row (when configured correctly) as its input argument. This allows the developer to precisely define how the values from that row should be channeled into and processed by the custom function. This approach dramatically simplifies the required syntax and is often cleaner and more maintainable compared to traditional iteration methods, such as explicit loops, which are generally slower and less Pythonic within the [Pandas](#) ecosystem.

To illustrate the core concept, assume we have previously defined a custom Python function named `f` that requires two distinct numerical inputs, perhaps corresponding to a player's `points` and `assists`. The basic syntax for applying this function row-wise and storing the output in a new column is structured as follows:

```
df = df.apply(lambda x: f(x.points, x.assists), axis=1)
```

In this powerful, single-line expression, the custom function `f` is executed using the specific values found in the **points** and **assists** columns of the DataFrame for every single row. The result of this calculation is efficiently collected and assigned to a newly created column, designated as **new_col**. This streamlined syntax is both efficient to execute and highly readable, solidifying its status as the preferred technique for performing row-wise, multi-column calculations.

Understanding the Crucial `axis` Parameter

Successful application of a function across multiple columns hinges entirely on the correct setting of the **axis** argument within the `apply()` method. The method inherently requires clear directionality: should the function be applied horizontally along the rows, or vertically down the columns?

When we explicitly set the parameter to **axis=1**, we are instructing [Pandas](#) to operate row-wise. This means that the input variable (conventionally named `x`) received by the [lambda function](#) represents each row of the DataFrame sequentially, presented as a Pandas Series object. Since `x` represents a single observation, we gain the crucial ability to access the individual column values within that row using simple dot notation, such as `x.points` or `x.assists`. This configuration is absolutely essential because the core requirement is a calculation where the inputs must originate from different features (columns) but belong to the same single observation (row).

Conversely, if the parameter is omitted or set to its default value, **axis=0**, the function would be applied column-wise. In the **axis=0** scenario, the variable `x` would represent an entire column (a Series), and the operation would be applied vertically down the column values. While this is highly suitable for aggregating data or transforming all values within a single feature (e.g., calculating the mean of a column), it is fundamentally incorrect for operations that aim to combine data across features on a record-by-record basis. Therefore, for any operation generating a new value per record based on multiple column inputs, the setting **axis=1** is mandatory.

Practical Example: Defining the DataFrame and Custom Function

To clearly demonstrate this powerful data processing technique, we will first establish a sample [DataFrame](#) containing simulated data for various basketball players. This illustrative example includes common statistical metrics such as the player's team, points scored, and assists made.

We begin by importing the necessary [Pandas](#) library and then constructing our foundational dataset:

```
import pandas as pd
```

```
#create DataFrame
```

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

```
'points': ,  
'assists': })  
  
#view DataFrame  
print(df)  
  
team points assists  
0 A 12 8  
1 A 15 10  
2 B 29 11  
3 B 22 11  
4 C 30 7  
5 C 41 14  
6 C 12 18
```

Next, suppose our analytical objective is to calculate a synthetic performance metric that effectively combines the player's offensive output. Specifically, we are required to perform the following two sequential mathematical operations for every player record:

First, calculate the product (multiplication) of the values found in the **points** and **assists** columns. Second, divide this resultant product by a standardized constant factor, which we will set to 3.

This bespoke logic must be clearly defined as a standard Python function using the **def** statement. This function, which we will simply name `f`, is designed to accept two arguments, corresponding directly to the two columns we intend to use in the calculation:

```
def f(first, second):  
    return (first*second) / 3
```

This function definition is intentionally simple yet highly robust. It efficiently accepts two inputs, performs the specified arithmetic calculation, and utilizes the **return** statement to provide the single, precise result needed for our new calculated performance column.

Implementing the Multi-Column Calculation

With the sample DataFrame successfully created and the custom function `f` explicitly defined, the final step involves integrating these components seamlessly using the [apply\(\)](#) method. We must carefully map the arguments required by our function (`first` and `second`) to the specific columns (`points` and `assists`) within the DataFrame structure.

We execute this crucial integration by assigning the output generated by the [apply\(\)](#) operation

directly to a new column, which we label `new_col`. The [lambda function](#) ensures that for every row processed (`x``), the corresponding column values `x.points` and `x.assists` are correctly extracted and passed as arguments to our custom function `f``.

```
#apply function to multiple columns in DataFrame  
df = df.apply(lambda x: f(x.points, x.assists), axis=1)
```

```
#view updated DataFrame  
print(df)
```

```
team points assists new_col  
0 A 12 8 32.000000  
1 A 15 10 50.000000  
2 B 29 11 106.333333  
3 B 22 11 80.666667  
4 C 30 7 70.000000  
5 C 41 14 191.333333  
6 C 12 18 72.000000
```

The resulting DataFrame now clearly includes the new feature, `new_col`, containing the calculated performance metric. We can verify the accuracy of these calculations by manually checking the first few rows against the defined mathematical logic:

Row 0 Calculation: (12 points multiplied by 8 assists) divided by 3 results in $96 / 3$, yielding **32.00**.

Row 1 Calculation: (15 points multiplied by 10 assists) divided by 3 results in $150 / 3$, yielding **50.00**.

Row 2 Calculation: (29 points multiplied by 11 assists) divided by 3 results in $319 / 3$, approximating **106.33**.

As demonstrated, the use of `axis=1` successfully ensured that the custom function was applied horizontally across the columns for each row, producing the precise, desired result in the newly appended feature column.

Performance Considerations and Advanced Usage

The newly generated `new_col` holds the output of our custom function, effectively reflecting a complex transformation derived from the input variables `points` and `assists`. This flexibility is the primary and most significant advantage of utilizing the `apply()` method for multi-column operations: it empowers the data analyst to define virtually any level of complexity in logic--ranging from simple arithmetic to highly sophisticated statistical formulas--and execute this logic efficiently across the entire dataset.

It is important to note that while our example used only two input columns, your custom Python function can be designed to accept any number of arguments, provided those arguments are correctly mapped to their corresponding columns within the [DataFrame](#) via the [lambda function](#). Furthermore, the body of the function itself can incorporate advanced structures, such as conditional statements (`if/else` logic) or calls to external Python libraries, enabling highly sophisticated, record-level decision-making processes during data preparation.

A crucial consideration, however, is performance. While the **`apply()`** method is highly intuitive and flexible, it is generally slower than optimized [vectorization](#) methods built directly into [Pandas](#) or NumPy (e.g., performing a calculation like `df * df`). For processing extremely large datasets where computational speed is paramount, data professionals should always verify if the required operation can be accomplished using native vectorized methods before defaulting to **`apply()`**, as the latter involves slower, Python-level iteration. Nevertheless, for ensuring code readability and implementing complex or non-vectorizable logic, **`apply()`** remains the most accessible and standard solution.

Note: For deeper technical insights, you can consult the complete documentation for the **`apply`** function available in the pandas official documentation.

Additional Resources for Pandas Operations

The following curated tutorials explain how to perform other common operations in [Pandas](#), helping you further master data manipulation in Python:

Featured Posts

[5 Statistical Biases to Avoid](#)

April 25, 2024

[5 Free Statistics Courses for Beginners](#)

April 19, 2024

[5 MIT Statistics Courses That Are Free](#)

April 18, 2024

[5 Free Books to Learn Statistics](#)

April 18, 2024

[How to Use the info\(\) Method in Pandas](#)

April 12, 2024

[How to Use pct_change\(\) in Pandas](#)

April 12, 2024