

Learning to Combine Data: A Guide to Appending Multiple Pandas DataFrames in Python

Authored by
Mohammed Iooti

October 29, 2025

RECOMMENDED CITATION

Mohammed Iooti (2025). *Learning to Combine Data: A Guide to Appending Multiple Pandas DataFrames in Python*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=5374>

In the realm of data science and analysis, the need to consolidate disparate datasets into a single, unified structure is constant. To efficiently combine multiple [Pandas DataFrames](#) (DFs) into a single, cohesive unit, a fundamental syntax leveraging the power of the [Pandas](#) library is utilized. This method is absolutely essential for complex data aggregation projects, allowing data professionals to seamlessly integrate various sources of information into one unified [Pandas DataFrame](#) for streamlined processing and analysis.

import pandas as pd

```
#append multiple DataFrames  
df_big = pd.concat(, ignore_index=True)
```

This powerful syntax employs the highly versatile [pd.concat\(\)](#) function to append **df1**, **df2**, and **df3** sequentially, resulting in a new, much larger [Pandas DataFrame](#) designated as **df_big**. The inclusion of the critical parameter `ignore_index=True` ensures that the output maintains a clean, continuous, and unambiguous [index](#), preventing potential conflicts that arise from duplicating row labels across different source DataFrames. Throughout this comprehensive guide, we will delve into the nuances of the [pd.concat\(\)](#) operation, exploring practical scenarios and key parameters that govern its behavior, such as index management and the alignment of disparate columns.

Core Functionality of the pd.concat() Method

[Pandas](#) is recognized as an indispensable library within the [Python](#) ecosystem, serving as the industry standard for robust data manipulation and analysis. Its primary data structure, the **DataFrame**, provides a sophisticated two-dimensional, mutable, and potentially heterogeneous tabular structure complete with labeled axes, making it ideal for handling complex datasets. A frequent requirement in data preprocessing involves combining multiple DataFrames, whether through vertical stacking (appending rows) or horizontal joining (adding columns), and [pd.concat\(\)](#) is the essential tool for managing these operations efficiently.

The [pd.concat\(\)](#) function acts as the cornerstone for merging [Pandas](#) objects along a specified axis. Crucially, it is distinct from database-style merge operations; `concat` is specifically engineered for stacking or "gluing together" objects either along the row axis (default) or the column axis. Its versatility stems from a range of optional parameters that grant fine-grained control over the concatenation outcome, including how indices are managed and how column names that do not perfectly align across all input DataFrames are treated. This flexibility ensures that the function can accommodate a wide array of data consolidation needs.

The fundamental method for utilizing [pd.concat\(\)](#) involves passing a list of [Pandas DataFrames](#) that are intended for combination. By default, the function performs a vertical concatenation (where

`axis=0`), meaning rows from the subsequent DataFrames are added to the end of the previous ones. The process attempts to align and preserve all columns found across the entire list of inputs. This behavior makes `concat` the perfect utility for consolidating datasets that share a similar schema or structure but contain distinct observations or records collected from different sources or time periods.

Vertical Stacking: A Practical Row-Wise Example

The most prevalent application of the `pd.concat()` function is the vertical appending of [Pandas DataFrames](#), which is achieved by stacking the rows of one DF directly beneath those of another. This operation is implicitly controlled by the `axis` parameter, which defaults to `0`, explicitly commanding concatenation along the row axis. When `axis=0` is active, [Pandas](#) aligns the DataFrames based on their column headers and then stacks them row-wise, ensuring structural consistency across the combined dataset.

Imagine a scenario in which business data is fragmented across several monthly files, each stored as a separate [Pandas DataFrame](#), yet all possessing an identical column layout. The objective is to merge these files into one single, comprehensive dataset for quarterly reporting. The following example demonstrates how to efficiently execute this common task by first establishing three distinct DataFrames and subsequently combining them into a unified structure using a single, optimized `concat` call.

import pandas as pd

```
#create three DataFrames
df1 = pd.DataFrame({'player': ,
'points':})

df2 = pd.DataFrame({'player': ,
'points':})

df3 = pd.DataFrame({'player': ,
'points':})

#append all DataFrames into one DataFrame
df_big = pd.concat(, ignore_index=True)

#view resulting DataFrame
print(df_big)

player points
0 A 12
```

1 B 5
2 C 13
3 D 17
4 E 27
5 F 24
6 G 26
7 H 27
8 I 27
9 J 12
10 K 9
11 L 5
12 M 5
13 N 13
14 O 17

As clearly demonstrated by the resulting output, the concatenation operation successfully merged the rows from **df1**, **df2**, and **df3** into a single, comprehensive [Pandas DataFrame](#) named **df_big**. The rows originating from each source DataFrame are stacked sequentially, creating one large structure that encapsulates all individual records, thereby providing a consolidated view of the entire dataset. The use of the argument `ignore_index=True` was pivotal here; it instructed [Pandas](#) to discard the original, potentially non-unique [index](#) numbers from the source DFs and instead generate a brand new, continuous [index](#) for the aggregated DataFrame, starting cleanly from zero. This is the recommended practice for simple data appending, as it ensures a clean and unambiguous sequential [index](#) across the combined result.

Managing Indices: The Importance of `ignore_index`

The [index](#) of a [Pandas DataFrame](#) is a critical structural element that facilitates efficient data retrieval and alignment. When concatenating multiple DFs, the method used to handle the index can directly influence the usability and integrity of the final dataset. The `pd.concat()` function addresses this concern via the `ignore_index` parameter.

If `ignore_index` is set to `False` (which is its default state), [Pandas](#) diligently preserves the original [index](#) values from every DataFrame in the concatenation list. While this preservation might be necessary in advanced data workflows where the index itself holds intrinsic semantic value (such as unique IDs or specific timestamps), it most frequently results in duplicate [index](#) values within the combined [Pandas DataFrame](#). This duplication can introduce significant ambiguity during subsequent operations, complicating data filtering, selection using `.loc`, and data alignment processes.

To demonstrate the potential pitfall of index duplication, let us consider an example where two DataFrames, **df1** and **df2**, are created with custom, overlapping [index](#) values. When these DFs are concatenated without explicitly setting `ignore_index=True`, the original indices are retained, resulting in the following structure where index labels are repeated across the combined dataset:

```
import pandas as pd
```

```
#create two DataFrames with indices
```

```
df1 = pd.DataFrame({'player': ,  
'points':},  
index=)
```

```
df2 = pd.DataFrame({'player': ,  
'points':},  
index=)
```

```
#stack the two DataFrames together
```

```
df_big = pd.concat()
```

```
#view resulting DataFrame
```

```
print(df_big)
```

```
player points
```

```
0 A 12
```

```
1 B 5
```

```
2 C 13
```

```
3 D 17
```

```
4 E 27
```

```
2 F 24
```

```
4 G 26
```

```
5 H 27
```

```
6 I 27
```

```
9 J 12
```

The resulting [Pandas DataFrame](#) clearly contains duplicate [index](#) values, such as 2 and 4, each appearing twice. Although computationally valid, this redundancy can introduce inconsistencies when attempting to select or manipulate data based on the index label. For the vast majority of simple data aggregation tasks where data origin tracking is unnecessary, it is strongly recommended that developers set `ignore_index=True`. This practice guarantees that the combined DataFrame possesses a fresh, sequential, and unambiguous index, thereby simplifying subsequent data manipulation steps and improving code reliability.

Handling Structural Differences with the join Parameter

When working with real-world data, it is rare for all DataFrames intended for concatenation to share an identical set of columns. Datasets collected from different systems might track slightly different metrics or use inconsistent naming conventions for attributes. To manage these column discrepancies effectively during the concatenation process, the `pd.concat()` function provides the essential `join` parameter.

The default setting for the `join` parameter is `'outer'`. When this setting is used, [Pandas](#) performs a union of all column names found across the input [Pandas DataFrames](#). Consequently, all columns from all source DFs will be present in the final combined DataFrame. For any row where a column did not exist in the original source, [Pandas](#) automatically inserts the missing value indicator, `NaN` (Not a Number). This outer join approach is critical because it ensures that absolutely no data is lost during the concatenation, making it the safest default choice, particularly during exploratory data analysis.

Alternatively, setting `join='inner'` instructs [Pandas](#) to calculate the intersection of all column names. In this scenario, only those columns that are explicitly common to *all* input DataFrames will be retained and included in the final concatenated result. Any columns that are unique to one or more of the source [Pandas DataFrames](#) are systematically discarded. This method is valuable when the analysis requires strict alignment across attributes, ensuring that the output only contains shared features, though it risks omitting potentially relevant, unique data points.

To clearly illustrate the fundamental operational difference between the `'outer'` and `'inner'` join types when concatenating [Pandas DataFrames](#) that possess non-identical columns, consider the following example involving separate sales and cost datasets:

```
import pandas as pd

# Create DataFrames with different columns
df_sales = pd.DataFrame({'product': ,
'revenue': })

df_costs = pd.DataFrame({'product': ,
'cost': })

# Concatenate with join='outer' (default)
df_outer = pd.concat(, ignore_index=True, join='outer')
print("Outer Join Result:")
print(df_outer)

print("\n")
```

```
# Concatenate with join='inner'
df_inner = pd.concat(, ignore_index=True, join='inner')
print("Inner Join Result:")
print(df_inner)

# Expected Output for Outer Join:
# product revenue cost
# 0 A 100.0 NaN
# 1 B 150.0 NaN
# 2 C 200.0 NaN
# 3 B NaN 50.0
# 4 C NaN 75.0
# 5 D NaN 120.0

# Expected Output for Inner Join:
# product
# 0 A
# 1 B
# 2 C
# 3 B
# 4 C
# 5 D
```

The `'outer'` join result successfully combines all columns (**product**, **revenue**, and **cost**) from both source DataFrames, utilizing [NaN](#) values to denote entries that were missing in the originals. Conversely, the `'inner'` join strictly retains only the **product** column, as it is the sole attribute common to both **df_sales** and **df_costs**, thereby discarding the unique revenue and cost metrics. A solid understanding of the implications of the `join` parameter is absolutely critical for accurately and reliably combining DataFrames with varying structural definitions.

Organizing Source Data with Hierarchical Keys

When performing concatenation across multiple [Pandas DataFrames](#), it is often essential to preserve clear metadata regarding the origin of each row. The `pd.concat()` function provides the specialized `keys` parameter specifically for this tracking purpose. By supplying a list of descriptive keys, [Pandas](#) automatically generates a [MultiIndex](#) (also known as a hierarchical [index](#)) within the resulting DataFrame. This MultiIndex incorporates an outer level corresponding to the provided keys, effectively tagging the source DataFrame for every block of appended rows.

The `keys` parameter proves particularly invaluable when combining datasets sourced from

disparate origins or collected during distinct time periods, necessitating a structured mechanism to quickly identify the original [Pandas DataFrame](#) for each record post-concatenation. This feature embeds an extra layer of organizational metadata directly into the DataFrame's [index](#), significantly simplifying subsequent tasks such as filtering based on source, grouping data by origin, or performing origin-specific analytical procedures.

The following example illustrates the application of the `keys` parameter to create a [MultiIndex](#). We will reuse a set of three DataFrames--**df1**, **df2**, and **df3**--and assign distinct, descriptive keys to each one during the concatenation process:

import pandas as pd

```
# Create three DataFrames
df1 = pd.DataFrame({'player': ,
'score': })

df2 = pd.DataFrame({'player': ,
'score': })

df3 = pd.DataFrame({'player': ,
'score': })

# Concatenate with keys to create a MultiIndex
df_multi_indexed = pd.concat(, keys=)

# View the resulting DataFrame with MultiIndex
print(df_multi_indexed)

# Expected Output:
# player score
# Group1 0 A 10
# 1 B 20
# 2 C 30
# Group2 0 D 25
# 1 E 35
# Group3 0 F 15
# 1 G 22
# 2 H 28
# 3 I 40
```

The output clearly displays the hierarchical [index](#), featuring 'Group1', 'Group2', and 'Group3' as the

distinct outer levels. Each group preserves the original index of its respective source [Pandas DataFrame](#) within the inner level. This sophisticated structure enables highly flexible data selection and aggregation based on both the original DataFrame's identity and its internal row index. When opting to use the `keys` parameter, it is generally prudent to omit `ignore_index=True`, as the resulting [MultiIndex](#) inherently resolves potential index conflicts by organizing them into distinct, identifiable tiers.

Recommended Practices and Alternatives

Although [pd.concat\(\)](#) stands as the most flexible and preferred method for appending [Pandas DataFrames](#), data practitioners should be aware of historical methods and established best practices. It is important to note that the older `DataFrame.append()` method, which provided similar row concatenation functionality, has been officially deprecated since [Pandas](#) version 1.4.0 and is scheduled for complete removal in future releases. The explicit recommendation from the [Pandas](#) development team is to exclusively use [pd.concat\(\)](#) due to its superior performance characteristics, enhanced flexibility, and consistent API design for combining all Pandas objects.

For scenarios involving exceptionally large datasets, performance optimization becomes a primary concern. While [pd.concat\(\)](#) is highly optimized, repeatedly concatenating hundreds or thousands of DataFrames sequentially within an iterative loop can still prove inefficient due to constant memory reallocation. A crucial performance best practice dictates that all DataFrames should first be collected and stored within a standard [Python](#) list. This list is then passed to [pd.concat\(\)](#) in a single, atomic call. This approach allows [Pandas](#) to optimize memory allocation and computation across the entire operation, dramatically increasing processing speed compared to repeated appending.

Finally, it is paramount to differentiate clearly between the processes of concatenation and merging. The [pd.concat\(\)](#) function is utilized for "stacking" or "gluing" DataFrames, predominantly operating row-wise to add records. Conversely, [Pandas](#) offers the separate `pd.merge()` function, which is designed for combining DataFrames based on matching values in common columns or indices, analogous to SQL join operations. The selection between `concat` and `merge` depends entirely on the analytical objective: use `concat` when you intend to add more rows (records) and use `merge` when you aim to add more columns (attributes) based on shared keys. For the specific task of simple row appending, [pd.concat\(\)](#) remains the most appropriate and powerful choice.

Conclusion

Mastering the ability to effectively combine multiple [Pandas DataFrames](#) is a foundational requirement for any data professional working within the [Python](#) data stack. The [pd.concat\(\)](#)

function provides a robust, flexible, and efficient mechanism for achieving this goal, granting users granular control over how their data is aggregated and structured.

By fully understanding and strategically applying its key operational parameters--including the `axis` for determining direction, `ignore_index` for ensuring a sequential and clean resulting `index`, `join` for managing disparate columns, and `keys` for generating a descriptive `Multindex`--you can ensure that all data concatenation operations are precise, predictable, and perfectly aligned with your intended analytical goals. Achieving this level of mastery will profoundly streamline data preparation workflows and enhance the overall reliability of your data analyses.

Additional Resources

[How to Add an Empty Column to a Pandas DataFrame](#)

[How to Insert a Column Into a Pandas DataFrame](#)

[How to Export a Pandas DataFrame to Excel](#)