

How to Easily Compare Two DataFrames Row by Row

Authored by
stats writer

November 21, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Compare Two DataFrames Row by Row*.
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=98827>

Comparing two DataFrames row by row is a common yet crucial task in data analysis, particularly when tracking changes, validating data pipelines, or performing quality assurance checks between two versions of a dataset. While manual looping is possible, modern data libraries provide highly optimized functions that streamline this process efficiently. The primary goal is often to identify every single mismatch between corresponding rows based on shared indices or positional alignment.

To successfully perform a row-by-row comparison, one must first ensure that both DataFrames adhere to certain structural prerequisites, such as having the same number of columns and consistent column names. Once alignment is verified, powerful comparison utilities--such as the compare() function provided by the pandas library--can be employed. Depending on whether you need to see only the differences or retain the entire shape of the input data, specific arguments within these functions, like equality operators or specialized numerical checks like np.allclose(), will dictate the final output.

Prerequisites for Effective DataFrame Comparison

Before initiating any structural comparison, analysts must confirm that the two DataFrames are ready for direct row-wise alignment. The most fundamental requirement is that both structures must possess identical dimensions or at least be alignable based on their row indices. If the indices are mismatched or unsynchronized, the comparison will fail to accurately pair corresponding observations, leading to incorrect difference identification. It is often necessary to reset or sort the indices of both DataFrames prior to comparison.

Furthermore, column names and data types must be consistent across both structures. A comparison operation relies on pairing column 'A' in the first DataFrame (df1) with column 'A' in the second DataFrame (df2). If a column name differs, or if the data types for corresponding columns are incompatible (e.g., comparing an integer column to a string column), the comparison might raise an error or produce misleading results. Ensuring data cleanliness and schema consistency is therefore a critical preparatory step.

The Modern Approach: Utilizing the pandas compare() Function

The pandas library offers the highly effective `.compare()` method, which is specifically designed for structural comparisons between two like-indexed objects. This function supersedes older, less efficient techniques involving merging, boolean masking, or manually written loops that iterate over the rows. The `.compare()` method provides a standardized and optimized way to highlight discrepancies, returning a resultant DataFrame that explicitly details where the differences lie.

The flexibility of the compare() function comes from its parameters, particularly `keep_equal` and `keep_shape`. By default, the function attempts to minimize the output by showing only the differing

values. However, using `keep_equal=True` allows the equal values to be retained in the output, which is helpful for context, especially when debugging. The parameter `align_axis=0` ensures that the comparison is performed row-wise, treating the rows of the two DataFrames as corresponding pairs, which is essential for accurate row-by-row analysis.

Setting Up the Demonstration DataFrames

To illustrate the practical application of the `.compare()` method, we will define two sample DataFrames, `df1` and `df2`, which share mostly common data but contain specific differences in rows 1 and 3. These examples demonstrate how the comparison function correctly isolates these specific changes.

import pandas as pd

```
#create first DataFrame
```

```
df1 = pd.DataFrame({'team': ,  
'points': ,  
'assists': })
```

```
print(df1)
```

```
team points assists
```

```
0 A 18 5  
1 B 22 7  
2 C 19 7  
3 D 14 9
```

```
#create second DataFrame
```

```
df2 = pd.DataFrame({'team': ,  
'points': ,  
'assists': })
```

```
print(df2)
```

```
team points assists
```

```
0 A 18 5  
1 B 30 7  
2 C 19 7  
3 E 20 9
```

A quick inspection reveals two key disparities: at index 1, the 'points' column differs (22 vs 30); and

at index 3, both the 'team' column ('D' vs 'E') and the 'points' column (14 vs 20) are different. The 'assists' column remains identical across all corresponding rows.

Example 1: Compare DataFrames and Only Keep Rows with Differences

In many analytical scenarios, the goal is to filter out all identical rows and focus exclusively on the points of divergence. This approach is highly effective for auditing changes or identifying anomalies between datasets. We achieve this by using the `compare()` function with `keep_equal=True`, which retains the equal values for contextual completeness, but without the `keep_shape=True` argument, which causes `pandas` to drop rows where no differences were found.

The following code executes the comparison, ensuring that only rows with at least one column mismatch are returned in the resulting `df_diff` `DataFrame`. The output format features a `pandas` `MultIndex` in the columns, distinguishing between values originating from the first `DataFrame` (labeled `self`) and the second (labeled `other`).

```
#compare DataFrames and only keep rows with differences  
df_diff = df1.compare(df2, keep_equal=True, align_axis=0)
```

```
#view results
```

```
print(df_diff)
```

```
team points
```

```
1 self B 22
```

```
other B 30
```

```
3 self D 14
```

```
other E 20
```

The resulting `DataFrame`, `df_diff`, clearly indicates that only rows indexed 1 and 3 contained differences. Row 0 and Row 2, being identical across all columns, were successfully excluded from the output. The primary effect of using `keep_equal=True` is demonstrated by the fact that the 'team' value in index 1 is shown as 'B' for both `self` and `other`, even though the 'points' value differs. If `keep_equal` had been set to `False`, the identical 'team' value would have been replaced by `NaN`.

Analyzing the Output of Difference Isolation

The structure of the comparison output is critical for interpretation. The columns are structured as a `MultIndex`, where the top level consists of the original column names (e.g., 'team', 'points'), and the secondary level denotes the source of the data: `self` refers to the values from `df1` (the object

calling the method), and `other` refers to values from `df2` (the object being compared against). This clear labeling ensures traceability of the discrepancies.

We observe that the 'assists' column is entirely missing from the output. This is a crucial feature of difference isolation: if an entire column contains no differences across the returned rows, pandas automatically excludes it to minimize redundancy and highlight only the relevant fields. Since 'assists' was the same in all rows of `df1` and `df2`, and since rows 0 and 2 were dropped (which were fully equal), the 'assists' column provided no difference information for the remaining rows and was thus dropped.

Let us examine the identified differences in detail:

The row at index position 1 is retrieved because of the divergence in the 'points' column:

The **self** row (`df1`) shows **22** points.

The **other** row (`df2`) shows **30** points.

The row at index position 3 is retrieved because of differences in both 'team' and 'points' columns:

The **self** row (`df1`) shows team **D** and **14** points.

The **other** row (`df2`) shows team **E** and **20** points.

It is important to understand that the argument `keep_equal=True` instructs pandas to maintain values that are equal within the scope of the differing rows. Without this, any equal values in those rows would be represented as Not a Number (NaN), making the output less informative regarding the original data content.

Example 2: Comparing DataFrames While Retaining Full Shape

In scenarios requiring complete data visualization or when the analyst must confirm the alignment of every row regardless of equality, it is necessary to retain the full shape of the original DataFrames. This method ensures that the resultant comparison structure possesses the same indices and number of rows as the input structures, with differences highlighted inline.

To achieve this, we introduce the `keep_shape=True` argument alongside `keep_equal=True` within the `compare()` function. This forces the inclusion of all rows, even those where `df1` and `df2` are perfectly identical across all columns. This is useful for debugging and visual verification of alignment.

#compare DataFrames and keep all rows

```
df_diff = df1.compare(df2, keep_equal=True, keep_shape=True, align_axis=0)
```

```
#view results
print(df_diff)

team points assists
0 self A 18 5
other A 18 5
1 self B 22 7
other B 30 7
2 self C 19 7
other C 19 7
3 self D 14 9
other E 20 9
```

The resulting `DataFrame` now contains all four index positions (0, 1, 2, and 3) from the original structures. Furthermore, because we retained the full shape, the 'assists' column is also included, as it is a part of the structural definition, even though its values are identical in the `self` and `other` rows. This comprehensive output facilitates a complete audit trail.

Understanding the Impact of Comparison Parameters

The key difference between the two methods lies in the behavior related to equality and output size. When `keep_shape=True` is used, the function guarantees that the resulting structure mirrors the full index range of the input structures. This is particularly important when intending to merge the difference results back into another large dataset or when performing visual inspection where missing rows could imply misalignment.

Conversely, when `keep_shape` is omitted (as in Example 1), the function acts as a filter, removing any row index where `df1` and `df2` are identical. This filtering mechanism is highly efficient for large datasets where differences are rare, drastically reducing memory usage and processing time by focusing only on the deviations.

Advanced Comparison Techniques

While the `compare()` function is the best tool for structural comparison, other techniques remain relevant for specific use cases. For instance, comparing two `DataFrames` simply for absolute equality often involves using boolean indexing, such as `(df1 == df2).all(axis=1)`, which returns a boolean series indicating row-level perfect matches.

However, when dealing with floating-point numbers, direct equality checks (`==`) are often unreliable due to precision issues inherent in computer arithmetic. In such cases, one should rely on NumPy

functions like `np.allclose()`, which checks if two arrays are element-wise equal within a specified tolerance (absolute tolerance and relative tolerance). This ensures that minor computational differences do not register as significant disparities, providing a more robust numerical comparison.

The methods detailed here provide robust and clear ways to compare two DataFrames row by row using the optimized functionalities available in the pandas library. By carefully managing the `keep_equal` and `keep_shape` parameters, analysts can tailor the comparison output to meet specific data validation or auditing requirements.

ARABPSYCHOLOGY.COM