

How to Fix You are trying to merge on object and int64 columns?

Authored by
stats writer

November 29, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Fix You are trying to merge on object and int64 columns?*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=101920>

Data integration is a cornerstone of modern data analysis. When working with complex datasets, analysts frequently need to combine information from multiple sources or tables. In the Python ecosystem, the **pandas** library provides powerful, flexible tools for this task, primarily through functions like **pandas.merge()**. The core goal of merging is to align rows from two different **DataFrame** objects based on shared key columns. However, this process relies heavily on the principle of data homogeneity: the values being matched must be structurally and typographically identical.

The challenge arises when the key columns intended for merging possess dissimilar **data type** specifications. Specifically, attempting to match a column designated as an **int64** (a 64-bit integer format common in numerical processing) with a column designated as an object (often representing strings or mixed types in **pandas**) will invariably lead to errors. Successfully combining these columns requires a crucial preparatory step: converting the values within the object column to conform to the numerical structure of the **int64** column. Only after this explicit type conversion can the **pandas.merge()** function successfully execute the join operation, yielding a coherent, single **DataFrame** or table.

During data manipulation within the **pandas** environment, a developer frequently encounters errors when data structures are not perfectly aligned. One of the most common and structurally informative errors related to joining operations is the following:

ValueError: You are trying to merge on int64 and object columns.

If you wish to proceed you should use pd.concat

This specific **ValueError** serves as a direct alert from **pandas**, indicating a foundational incompatibility in the columns designated as keys for the merge operation. The error explicitly states that one key column is stored as an **int64** data type, typically used for whole numbers, while the corresponding key column in the second **DataFrame** is identified as an object, which is **pandas**' generic type often used for strings or heterogeneous data. Because relational joins require exact matches on the key, the library defaults to raising an exception rather than attempting an unreliable implicit conversion, which might lead to unexpected results or data loss. Understanding this distinction is the first step toward resolution.

Understanding the Core Problem: Data Type Mismatch

The root cause of the merge failure lies in how programming environments, particularly Python and **pandas**, handle and compare different **data type** representations. An **int64** value, such as the number 2021, is stored internally as a fixed-width binary representation. Conversely, the string '2021', which is classified as an object type, is stored as a sequence of characters, potentially requiring variable width. Even though they look identical to the human eye, the computer interprets

them as fundamentally different entities. When **pandas.merge()** attempts to match these keys, it performs a strict comparison based on the underlying structure.

When the data types differ, **pandas** is unable to guarantee that the comparison logic will work correctly across the entire dataset. For instance, comparing numerical integers is highly optimized, while comparing strings (objects) involves character-by-character matching and potentially different hashing algorithms. To maintain computational integrity and speed, **pandas** enforces the rule that the join keys must possess identical **data type** specifications across both source **DataFrame** objects. This necessity for homogeneity ensures reliable, predictable joins and prevents subtle bugs that could arise from silent, inaccurate type casting.

This **ValueError** is a safeguard built into the library. Rather than attempting a potentially lossy or slow conversion behind the scenes, **pandas** forces the user to make an explicit choice about data structure. This design choice prevents common errors where numerical keys might fail to match string keys due to minor formatting differences (e.g., '100' vs ' 100') that are invisible to the merge operation but critical to the integrity of the join.

Setting Up the Scenario: Reproducing the Error

To fully grasp the issue and its solution, it is helpful to establish a working example where this merge incompatibility manifests. We will create two distinct **DataFrame** structures designed to hold related information--sales and refunds--but crucially, with a deliberate mismatch in the format of their shared key column, 'year'. This scenario models a common occurrence in real-world data science, where data sourced from different systems (e.g., one from an SQL database, one from a CSV file) often results in inconsistent **data type** assignments.

Our first **DataFrame**, `df1`, will store the annual sales figures, with the 'year' column naturally interpreted as a numerical integer. Our second **DataFrame**, `df2`, will store refunds, but we will explicitly define the 'year' column as strings (objects) by wrapping the values in quotes. This subtle difference is the exact mechanism that triggers the subsequent error when we attempt the merge operation. The setup phase is vital for debugging, as knowing the initial state of the data is paramount.

How to Reproduce the Error

Observe the creation and structure of the two sample datasets below. Note the difference in the input types for the 'year' column between the definitions of `df1` (using integers) and `df2` (using strings/objects) in the initialization code.

```
import pandas as pd
```

```
#create DataFrame
df1 = pd.DataFrame({'year': ,
'sales': })
```

```
df2 = pd.DataFrame({'year': ,
'refunds': })
```

```
#view DataFrames
```

```
print(df1)
```

```
year sales
0 2015 500
1 2016 534
2 2017 564
3 2018 671
4 2019 700
5 2020 840
6 2021 810
```

```
print(df2)
```

```
year refunds
0 2015 31
1 2016 36
2 2017 40
3 2018 40
4 2019 43
5 2020 70
6 2021 62
```

Upon reviewing the data types (e.g., using `df1.dtypes`), we would confirm that the 'year' column in `df1` is identified as **int64**. Conversely, reviewing the types for `df2` would reveal the 'year' column is labeled as `object`. This critical difference immediately prohibits the standard merge operation. When the command to merge is executed using the incompatible key, the system fails gracefully, preventing potentially disastrous data corruption by issuing the explicit **ValueError**.

Now suppose we attempt to merge the two **DataFrame** objects using the mismatched key, expecting **pandas** to handle the join automatically:

```
#attempt to merge two DataFrames
```

```
big_df = df1.merge(df2, on='year', how='left')
```

ValueError: You are trying to merge on int64 and object columns.

If you wish to proceed you should use `pd.concat`

We receive a definitive **ValueError**. This error clearly signals that the `year` variable in the first **DataFrame** (`df1`) is an **int64**, while the corresponding `year` variable in the second **DataFrame** (`df2`) is designated as an object. This explicit failure mechanism is designed to force the user to address the **data type** conflict before proceeding with the integration.

Strategies for Data Type Harmonization

When confronted with a **data type** mismatch during a merge, the most robust solution is to explicitly convert one of the key columns so that its type matches the other. In this specific scenario, where the 'year' column in `df2` is an object (string) representing numerical data, the most logical and efficient approach is to cast this string representation into a numerical **int64** format. This conversion should always be performed on the column that holds the less restrictive **data type** (usually the object or string column), assuming the data contained within it is purely numerical.

The conversion process is critical: before applying the cast, the analyst must verify that all values within the object column can be safely interpreted as integers. If the column contains non-numeric characters, such as 'N/A' or mixed text and numbers, the standard `.astype(int)` operation will itself raise a **ValueError** or similar exception. Therefore, ensuring data cleanliness--handling missing values or non-conformant entries--is a prerequisite to successful type casting. For our current example, since the strings are pure numeric representations of years, the conversion is straightforward and reliable.

The flexibility of **pandas** allows us to use the `.astype()` method directly on the targeted column (a **Series** object). This method is designed to change the underlying **data type** of the entire Series in place or return a new series with the desired type. By specifying `.astype(int)`, we instruct **pandas** to interpret the string contents of the 'year' column in `df2` as numerical integers, aligning it perfectly with the **int64** type found in `df1`. This prepares the data for a successful, high-performance join.

Implementing the Fix and Successful Merge

The solution involves a single, targeted line of code executed before the **pandas.merge()** call. We specifically target the 'year' column in the second **DataFrame** (`df2`) and apply the type conversion function. This action immediately resolves the structural inconsistency that blocked the previous merge attempt, making the two key columns functionally equivalent for relational joining.

How to Fix the Error

The easiest way to fix this error is to simply convert the `year` variable in the second **DataFrame** to an integer using the `.astype(int)` method, and then perform the merge.

The following syntax demonstrates the conversion and the subsequent successful merge operation:

```
#convert year variable in df2 to integer
```

```
df2=df2.astype(int)
```

```
#merge two DataFrames
```

```
big_df = df1.merge(df2, on='year', how='left')
```

```
#view merged DataFrame
```

```
big_df
```

```
year sales refunds
```

```
0 2015 500 31
```

```
1 2016 534 36
```

```
2 2017 564 40
```

```
3 2018 671 40
```

```
4 2019 700 43
```

```
5 2020 840 70
```

```
6 2021 810 62
```

Immediately after the type conversion, the execution of the **pandas.merge()** function proceeds without error. Notice that we do not receive any **ValueError** this time, and we are successfully able to integrate the data from the two source tables into a single, cohesive **DataFrame** named `big_df`. The success of this operation underscores the fundamental rule of data integration in **pandas**: matching data types on key columns is non-negotiable for relational joins.

Alternative Scenarios and Considerations

While converting the object column to an **int64** is the canonical fix when dealing with numerical keys, it is important to consider alternative scenarios. If the object column contained meaningful strings that were *not* simply numerical representations (e.g., product IDs like 'P-1001' vs. numerical IDs like 1001), converting the string column to an integer would be impossible and incorrect. In such cases, the alternative fix would be converting the **int64** column in `df1` to an object (string) type using `.astype(str)`. The choice of which column to convert depends entirely

on the domain knowledge and the intended final **data type** for the merged result.

Another related issue occurs when the object column contains truly mixed data types, such as integers, floats, and strings that cannot be cleanly converted to a single numerical format. If we attempt `.astype(int)` on such a column, **pandas** will fail. In these complex situations, analysts often resort to more sophisticated data cleaning routines, such as using `pd.to_numeric(errors='coerce')` to convert valid numeric entries while replacing invalid entries with `NaN`, followed by handling the resulting missing values before the merge. This robust approach handles data imperfections gracefully, ensuring that only valid keys are used for the join.

The error message itself suggests an alternative path: "If you wish to proceed you should use `pd.concat`." It is crucial to understand that `pd.concat` performs concatenation (stacking rows or columns) rather than relational merging (joining based on keys). While `pd.concat` is often type-agnostic, using it to combine datasets that require alignment based on keys will generally lead to incorrect data structure unless the indices of both **DataFrame** objects are already perfectly aligned. Therefore, for robust, key-based integration, explicit type harmonization remains the superior methodology over blindly utilizing `pd.concat` as a workaround.

Summary and Best Practices for Data Quality

The **ValueError** encountered during the merge of **int64** and object columns is a highly instructive moment in data preparation. It reinforces the necessity of proactive data quality checks. Before any major transformation or integration step, analysts should routinely inspect the data types of all columns involved, especially those designated as merge keys. This can be quickly achieved using the `.dtypes` attribute on any **DataFrame**.

Implementing a strict data validation phase early in the analysis pipeline saves significant time downstream. Best practices dictate that columns intended to serve as unique identifiers or relational keys should be standardized to the same **data type**, preferably at the point of data ingestion. For numerical identifiers, the **int64** type is generally appropriate, ensuring efficient storage and comparison. If keys contain non-standard characters, explicit string (object) manipulation should be used to ensure uniformity.

In conclusion, while **pandas** is flexible, its merging functionality demands precision. When faced with the object/int64 mismatch, the solution is clear: use the `.astype()` method to enforce **data type** consistency on the key columns. This ensures that the **pandas.merge()** function can perform its duties reliably, producing a clean, integrated dataset ready for advanced statistical analysis or machine learning pipelines.