# Relations between data frames

**Definition of Data Frames**

A data frame is a two-dimensional, size-mutable, and heterogeneous tabular data structure with labeled axes (rows and columns). It is akin to a table in a relational database or an Excel spreadsheet. Each column in a data frame can hold different types of data (numeric, string, boolean, etc.), and the rows represent records or observations.

**Key Characteristics:**

- **Two-Dimensional**: Data frames have rows and columns.

- **Labeled Axes**: Both rows and columns have labels (names) which allow for easy identification and access.

- **Heterogeneous Data Types**: Different columns can contain different types of data.

- **Mutable Size**: Data frames can be dynamically resized by adding or removing rows and columns.
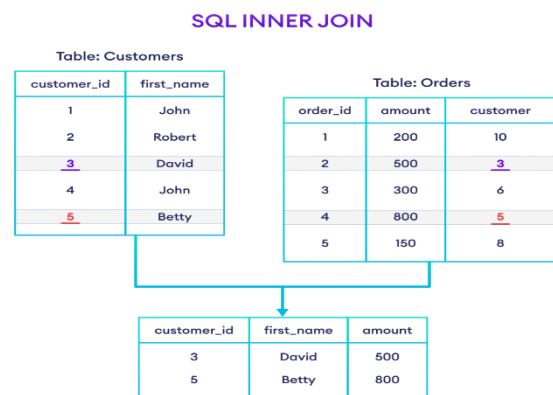
**Joining and Merging Data Frames**

Joining and merging data frames are essential operations in data science and machine learning, allowing for the combination of different datasets based on common keys. These operations are crucial for integrating data from multiple sources, enriching datasets, and performing comprehensive analyses.
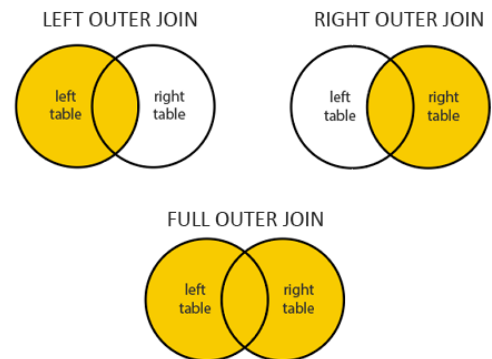
*Types of Joins*

1. **Inner Join**

- **Description**: Combines rows from two data frames where the join key exists in both data frames.
- **Result**: Only the rows with matching keys in both data frames are included.
- **Example**: If you have two data frames with customer data and order data, an inner join would return only the customers who have placed orders.

**SQL INNER JOIN**

Table: Customers

| customer_id | first_name |
|---|---|
| 1 | John |
| 2 | Robert |
| 3 | David |
| 4 | John |
| 5 | Betty |

Table: Orders

| order_id | amount | customer |
|---|---|---|
| 1 | 200 | 10 |
| 2 | 500 | 3 |
| 3 | 300 | 6 |
| 4 | 800 | 5 |
| 5 | 150 | 8 |

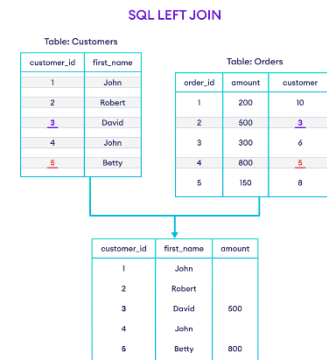| customer_id | first_name | amount |
|---|---|---|
| 3 | David | 500 |
| 5 | Betty | 800 |

## 2. Outer Join

- **Description**: Combines rows from two data frames, returning all rows from both data frames, with NaN in places where the join key is missing.

- **Result**: Includes all rows from both data frames, regardless of whether there is a match.

- **Example**: This join would return all customers and all orders, filling in NaN where a customer has no orders or an order has no corresponding customer.



## 3. Left Join

- **Description**: Combines rows from two data frames, returning all rows from the left data frame and the matched rows from the right data frame.
- **Result**: All rows from the left data frame are included, with NaN in the columns from the right data frame where there is no match.
- **Example**: Returns all customers and their orders, if any. Customers without orders will have NaN for order details.



## 4. Right Join

- **Description**: Combines rows from two data frames, returning all rows from the right data frame and the matched rows from the left data frame.
- **Result**: All rows from the right data frame are included, with NaN in the columns from the left data frame where there is no match.
- **Example**: Returns all orders and their corresponding customers, if any. Orders without corresponding customers will have NaN for customer details.

**Concatenating DataFrame**

In order to concat dataframe, we use concat() function which helps in concatenating a dataframe. We can concat a dataframe in many different ways, they are:

- Concatenating DataFrame using .concat()

**Concatenating DataFrame using .concat() :**

```
frames = [df, df1]

res1 = pd.concat(frames)
res1
```

**Output :**

As shown in the output image, we have created two dataframe after concatenating we get one dataframe.

| | Name | Age | Address | Qualification |
|---|---|---|---|---|
| 0 | Jai | 27 | Nagpur | Msc |
| 1 | Princi | 24 | Kanpur | MA |
| 2 | Gaurav | 22 | Allahabad | MCA |
| 3 | Anuj | 32 | Kannuaj | Phd |

| | Name | Age | Address | Qualification |
|---|---|---|---|---|
| 4 | Abhi | 17 | Nagpur | Btech |
| 5 | Ayushi | 14 | Kanpur | B.A |
| 6 | Dhiraj | 12 | Allahabad | Bcom |
| 7 | Hitesh | 52 | Kannuaj | B.hons |

| | Name | Age | Address | Qualification |
|---|---|---|---|---|
| 0 | Jai | 27 | Nagpur | Msc |
| 1 | Princi | 24 | Kanpur | MA |
| 2 | Gaurav | 22 | Allahabad | MCA |
| 3 | Anuj | 32 | Kannuaj | Phd |
| 4 | Abhi | 17 | Nagpur | Btech |
| 5 | Ayushi | 14 | Kanpur | B.A |
| 6 | Dhiraj | 12 | Allahabad | Bcom |
| 7 | Hitesh | 52 | Kannuaj | B.hons |

In Pandas, both join and merge functions are used to combine data from multiple data frames. However, they are used in slightly different contexts and offer different functionalities. Understanding the distinctions between join and merge can help you choose the appropriate method for combining your datasets. Here's a detailed comparison:

**merge Function**

The merge function in Pandas is used to combine two data frames based on one or more keys. It is similar to SQL joins and provides a high level of flexibility in terms of specifying how the join should be performed.

**Key Characteristics of merge:**

1. **Flexibility**: merge allows for various types of joins such as inner, outer, left, and right.

2. **Key Specification**: You can specify the columns to join on using the on, left_on, and right_on parameters.

3. **Suffixes**: It automatically handles overlapping column names by adding suffixes.

4. **Merge on Index**: Can merge on indexes using the left_index and right_index parameters.

**join Function**

The join function in Pandas is used to combine two data frames based on their indexes. It is particularly useful when you need to join along the index rather than specific columns.

**Key Characteristics of join:**

1. **Index-Based**: Primarily used for joining on indexes.
2. **Less Flexible**: Does not allow specification of multiple columns for joining like merge does.
3. **Simpler Syntax**: Easier to use when you are working with index-based joins.

In Pandas, join and merge are both used to combine data from multiple data frames, but they are not exactly the same. They serve different purposes and offer different functionalities, even though they can sometimes be used interchangeably. Here's a detailed comparison to highlight their differences and similarities:

**merge Function**

The merge function is very flexible and allows you to combine data frames based on one or more keys. It is analogous to SQL joins and is more general-purpose.

**Key Characteristics of merge:**

1. **Flexibility in Keys**: Can specify the columns to join on using the on, left_on, and right_on parameters.
2. **Join Types**: Supports various types of joins (inner, outer, left, right).
3. **Suffixes for Overlapping Columns**: Handles overlapping column names by adding suffixes (suffixes parameter).
4. **Join on Indexes**: Can merge on indexes using the left_index and right_index parameters.

**Example:**

python

Copy code

```
import pandas as pd


df1 = pd.DataFrame({

    'CustomerID': [1, 2, 3],

    'CustomerName': ['John', 'Jane', 'Joe']

})


df2 = pd.DataFrame({

    'CustomerID': [2, 3, 4],

    'OrderID': [101, 102, 103]

})


# Inner join on CustomerID

result = pd.merge(df1, df2, on='CustomerID', how='inner')

print(result)
```

**join Function**

The join function is primarily used to join data frames based on their indexes. It is convenient when you need to combine data frames using the index without explicitly specifying a key.

**Key Characteristics of join:**

1. **Index-Based Join**: Primarily used for joining on indexes.
2. **Simpler Syntax for Index Joins**: Easier to use when dealing with index-based joins.
3. **Join Types**: Also supports various types of joins (inner, outer, left, right), but these are applied to index joins.

4. **Suffixes for Overlapping Columns**: Handles overlapping column names using suffixes (lsuffix and rsuffix parameters).

**Key Differences**

1. **Join Keys**:
   - merge: Can join on specific columns or on indexes.
   - join: Primarily joins on indexes but can join on a key column by setting the on parameter (in DataFrame.join method).
2. **Syntax and Use Cases**:
   - merge: More flexible and general-purpose. Suitable for a wide range of use cases, especially when columns are the join keys.
   - join: Simpler and more convenient for index-based joins.
3. **Performance**:
   - Both merge and join are efficient, but merge offers more flexibility at the cost of slightly more complex syntax.
4. **Overlap Handling**:
   - Both handle overlapping columns by allowing suffixes to be added to the overlapping column names.

**When to Use Which**

- **Use merge**:
  - When you need to join on specific columns.
  - When you require the flexibility to specify multiple keys.
  - For SQL-like joins where specific join conditions are necessary.
- **Use join**:
  - When you need to join based on indexes.
  - For simpler joining requirements.
  - When working with data frames that are already indexed appropriately for the join.

In summary, data frames are a fundamental tool in data science and machine learning, providing a versatile and powerful structure for organizing, manipulating, and analyzing data. Their ability to handle diverse data types and support complex operations makes them indispensable for efficient and effective data analysis.