

Merge DataFrame or named Series objects with a database-style join

pandas.merge(left, right, how='inner', on=None, left_on=None, right_on=None, left_index=False,
right_index=False, sort=False, suffixes=('_x', '_y'), copy=True, indicator=False, validate=None)

- left: A DataFrame or named Series object.
- right: Another DataFrame or named Series object.
- on: Column or index level names to join on. Must be found in both the left and right DataFrame and/or Series objects. If not passed and left_index and right_index are False, the intersection of the columns in the DataFrames and/or Series will be inferred to be the join keys.
- **left_on**: Columns or index levels from the left DataFrame or Series to use as keys. Can either be column names, index level names, or arrays with length equal to the length of the DataFrame or Series.
- **right_on**: Columns or index levels from the right DataFrame or Series to use as keys. Can either be column names, index level names, or arrays with length equal to the length of the DataFrame or Series.
- **left_index**: If True, use the index (row labels) from the left DataFrame or Series as its join key(s). In the case of a DataFrame or Series with a MultiIndex (hierarchical), the number of levels must match the number of join keys from the right DataFrame or Series.
- right index: Same usage as left index for the right DataFrame or Series
- how: One of 'left', 'right', 'outer', 'inner', 'cross'. Defaults to inner. See below for more detailed description of each method.



pandas.merge(left, right, how='inner', on=None, left_on=None, right_on=None, left_index=False,
right_index=False, sort=False, suffixes=('_x', '_y'), copy=True, indicator=False, validate=None)

sort: Sort the result DataFrame by the join keys in lexicographical order. Defaults to True, setting to False will improve performance substantially in many cases.

suffixes: A tuple of string suffixes to apply to overlapping columns. Defaults to ('_x', '_y').

copy: Always copy data (default True) from the passed DataFrame or named Series objects, even when reindexing is not necessary. Cannot be avoided in many cases but may improve performance / memory usage. The cases where copying can be avoided are somewhat pathological but this option is provided nonetheless.

indicator: Add a column to the output DataFrame called _merge with information on the source of each row. _merge is Categorical-type and takes on a value of left_only for observations whose merge key only appears in 'left' DataFrame or Series, right_only for observations whose merge key only appears in 'right' DataFrame or Series, and both if the observation's merge key is found in both.

validate: string, default None. If specified, checks if merge is of specified type.

"one to one" or "1:1": checks if merge keys are unique in both left and right datasets.

[&]quot;one_to_many" or "1:m": checks if merge keys are unique in left dataset.

[&]quot;many to one" or "m:1": checks if merge keys are unique in right dataset.

[&]quot;many_to_many" or "m:m": allowed, but does not result in checks.



how options and SQL equivalent names

| Merge method | SQL Join Name | Description |
|--------------|------------------|-----------------------------------------------------|
| left | LEFT OUTER JOIN | Use keys from left frame only |
| right | RIGHT OUTER JOIN | Use keys from right frame only |
| outer | FULL OUTER JOIN | Use union of keys from both frames |
| inner | INNER JOIN | Use intersection of keys from both frames |
| cross | CROSS JOIN | Create the cartesian product of rows of both frames |



result = pd.merge(left, right, how="left", on=["key1", "key2"])

left

right

Result

| | keyl | key2 | Α | В |
|---|------|------|----|----|
| 0 | KO | KD | A0 | B0 |
| 1 | KO | K1 | A1 | B1 |
| 2 | K1 | KD | A2 | B2 |
| 3 | K2 | K1 | A3 | В3 |

|] | | keyl | key2 | С | D |
|---|---|------|------|----|----|
| | 0 | KD | KD | œ | D0 |
| | 1 | K1 | KO | C1 | D1 |
| I | 2 | K1 | KO | C2 | D2 |
| | 3 | K2 | KO | СЗ | D3 |

| | keyl | key2 | Α | В | С | D |
|---|------|------|----|----|-----|-----|
| 0 | KD | K0 | A0 | B0 | co | D0 |
| 1 | KD | K1 | A1 | B1 | NaN | NaN |
| 2 | K1 | K0 | A2 | B2 | C1 | D1 |
| 3 | K1 | KD | A2 | B2 | C2 | D2 |
| 4 | K2 | K1 | A3 | В3 | NaN | NaN |

result = pd.merge(left, right, how="right", on=["key1", "key2"])

left

right

Result

| | keyl | key2 | Α | В |
|---|------|------|----|----|
| 0 | KO | KD | A0 | B0 |
| 1 | KD | K1 | Al | B1 |
| 2 | K1 | KD | A2 | B2 |
| 3 | K2 | K1 | A3 | В3 |

| | | keyl | key2 | С | D | |
|---|---|------|------|----|----|--|
| 1 | 0 | K0 | KO | œ | D0 | |
| | 1 | K1 | KO | C1 | D1 | |
| | 2 | K1 | KO | C2 | D2 | |
| | 3 | K2 | KO | СЗ | D3 | |
| | | | | | | |

| | keyl | key2 | Α | В | С | D |
|---|------|------|-----|-----|----|----|
| 0 | KO | K0 | A0 | B0 | co | D0 |
| 1 | K1 | K0 | A2 | B2 | CI | D1 |
| 2 | K1 | KO | A2 | B2 | C2 | D2 |
| 3 | K2 | K0 | NaN | NaN | СЗ | D3 |

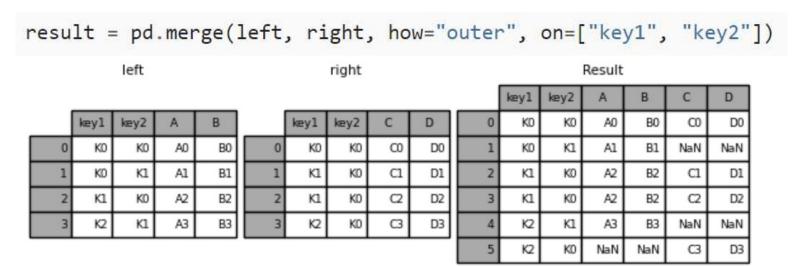


Pandas – merge() – example: left / right

Merge Dataframes

| In [9]: | | df_merge_left=pd.merge(df1,df3, how='left', on='id') df_merge_left | | | | | | | | | | |
|----------|----------|--------------------------------------------------------------------|------------|--------|------------------------|------|----------------|-------|---------|----------|-----|--|
| Out[9]: | | name_x | date | height | weight | age | name_y | hours | status | children | sex | |
| | id | | | | | | | | | | | |
| | 1373913 | Marisa Martins | 2013-02-05 | 155 | 48 | 45 | Marisa Martins | 3.0 | married | 2.0 | F | |
| | 1109818 | Rita Fonseca | 2018-08-28 | 166 | 54 | 45 | NaN | NaN | NaN | NaN | NaN | |
| [n [10]: | df_merge | e_right=pd.me e_right | erge(df1,d | f3, ho | w <mark>='rig</mark> h | it', | on='id') | | | | | |
| Out[10]: | | name_x | date | height | weight | age | name_y | hours | status | children | sex | |
| | id | | | | | | | | | | | |
| | 1767703 | NaN | NaT | NaN | NaN | NaN | Manuel Martins | 6 | single | 0 | М | |
| | 1373913 | Marisa Martins | 2013-02-05 | 155.0 | 48.0 | 45.0 | Marisa Martins | 3 | married | 2 | F | |
| | 1158813 | NaN | NaT | NaN | NaN | | Joana Freitas | 3 | widow | 1 | F | |





result = pd.merge(left, right, how="inner", on=["key1", "key2"])

left

| | keyl | keyl key2 A | | | |
|---|------|-------------|----|----|--|
| 0 | KO | KD | A0 | B0 | |
| 1 | KO | K1 | Al | B1 | |
| 2 | K1 | K0 | A2 | B2 | |
| 3 | K2 | K1 | A3 | В3 | |

right

| | keyl | key2 | С | D |
|---|------|------|----|----|
| 0 | K0 | KO | 00 | D0 |
| 1 | K1 | KD | Cl | D1 |
| 2 | K1 | KD | C2 | D2 |
| 3 | K2 | KO | СЗ | D3 |

Result

| | keyl | key2 | Α | В | О | D |
|---|------|------|----|----|----|----|
| 0 | KO | K0 | A0 | В0 | 00 | D0 |
| 1 | K1 | K0 | A2 | B2 | Cl | D1 |
| 2 | K1 | K0 | A2 | B2 | C2 | D2 |



Pandas – merge() – example: inner / outer

```
In [11]:
          df merge inner=pd.merge(df1,df3, how='inner',on='id')
          df merge inner
Out[11]:
                                       date height weight age
                                                                                     status children sex
                         name x
                                                                     name y
                                                                             hours
                 id
           1373913 Marisa Martins 2013-02-05
                                               155
                                                       48
                                                            45 Marisa Martins
                                                                                 3 married
                                                                                                  2
                                                                                                       F
          df_merge_outer=pd.merge(df1,df3, how='outer', on='id')
In [12]:
          df merge outer
Out[12]:
                                           height weight age
                                                                              hours
                         name x
                                                                      name y
                                                                                      status children
                 id
           1373913
                    Marisa Martins 2013-02-05
                                                                 Marisa Martins
                                                                                 3.0 married
                                                                                                 2.0
                                                                                                        F
                                             155.0
                                                      48.0 45.0
           1109818
                     Rita Fonseca 2018-08-28
                                                      54.0 45.0
                                                                                                NaN
                                             166.0
                                                                         NaN
                                                                                NaN
                                                                                        NaN
                                                                                                     NaN
           1767703
                                                                Manuel Martins
                            NaN
                                                      NaN NaN
                                                                                 6.0
                                                                                       single
                                                                                                 0.0
                                                                                                        M
                                       NaT
                                              NaN
           1158813
                            NaN
                                       NaT
                                              NaN
                                                      NaN
                                                           NaN
                                                                  Joana Freitas
                                                                                 3.0
                                                                                      widow
                                                                                                 1.0
                                                                                                        F
```



result = pd.merge(left, right, how="cross")

left right Result

| | keyl | key2 | Α | В | | keyl | key2 | С | D |
|---|------|------|----|----|---|------|------|----|----|
| 0 | KD | KD | A0 | B0 | 0 | K0 | K0 | co | D0 |
| 1 | KD | K1 | A1 | B1 | 1 | K1 | KD | C1 | D1 |
| 2 | K1 | KO | A2 | B2 | 2 | K1 | KD | C2 | D2 |
| 3 | K2 | K1 | A3 | В3 | 3 | K2 | KO | СЗ | D3 |

| | key1_x | key2_x | А | В | key1_y | key2_y | С | D |
|----|--------|--------|----|-----|--------|--------|---|----|
| 0 | KD | KD | AD | B0 | KD | KD | В | D0 |
| 1 | KD | KD | AD | BO | кі | KD | а | D1 |
| 2 | KD | KD | AD | B0 | кі | KD | Ŋ | D2 |
| 3 | KD | KD | AD | В0 | K2 | KD | Ω | D3 |
| 4 | KD | кі | A1 | B1. | KD | KD | В | D0 |
| 5 | KD | кі | A1 | B1 | кі | KD | а | D1 |
| 6 | KD | кі | A1 | B1 | кі | KD | Q | D2 |
| 7 | KD | кі | A1 | B1. | K2 | KD | В | D3 |
| 8 | кі | KD | A2 | B2 | KD | KD | 8 | D0 |
| 9 | кі | KD | A2 | B2 | и | KD | а | D1 |
| 10 | КІ | KD | A2 | B2 | кт | KD | a | D2 |
| 11 | КІ | KD | A2 | B2 | K2 | KD | В | D3 |
| 12 | K2 | кт | A3 | B3 | KD | KD | 8 | D0 |
| 13 | K2 | кі | A3 | В3 | кі | KD | п | D1 |
| 14 | K2 | кі | A3 | В3 | кі | Ю | a | D2 |
| 15 | K2 | KI | А3 | B3 | K2 | KD | З | D3 |



Pandas – merge() – example cross

```
In [13]: df_merge_cross = pd.merge(df1, df3, how='cross')
    df_merge_cross
```

Out[13]:

| | name_x | date | height | weight | age | name_y | hours | status | children | sex |
|---|----------------|------------|--------|--------|-----|----------------|-------|---------|----------|-----|
| 0 | Marisa Martins | 2013-02-05 | 155 | 48 | 45 | Manuel Martins | 6 | single | 0 | М |
| 1 | Marisa Martins | 2013-02-05 | 155 | 48 | 45 | Marisa Martins | 3 | married | 2 | F |
| 2 | Marisa Martins | 2013-02-05 | 155 | 48 | 45 | Joana Freitas | 3 | widow | 1 | F |
| 3 | Rita Fonseca | 2018-08-28 | 166 | 54 | 45 | Manuel Martins | 6 | single | 0 | М |
| 4 | Rita Fonseca | 2018-08-28 | 166 | 54 | 45 | Marisa Martins | 3 | married | 2 | F |
| 5 | Rita Fonseca | 2018-08-28 | 166 | 54 | 45 | Joana Freitas | 3 | widow | 1 | F |



Pandas – join()

Join columns with other DataFrame either on index or on a key column

DataFrame.join(other, on=None, how='left', lsuffix=", rsuffix=", sort=False)

other: DataFrame, Series, or list of DataFrame

Index should be similar to one of the columns in this one. If a Series is passed, its name attribute must be set, and that will be used as the column name in the resulting joined DataFrame.

on: str, list of str, or array-like, optional

Column or index level name(s) in the caller to join on the index in other, otherwise joins index-on-index. If multiple values given, the other DataFrame must have a MultiIndex. Can pass an array as the join key if it is not already contained in the calling DataFrame. Like an Excel VLOOKUP operation.

how: {'left', 'right', 'outer', 'inner'}, default 'left'

How to handle the operation of the two objects.

- left: use calling frame's index (or column if on is specified)
- right: use other's index.
- outer: form union of calling frame's index (or column if on is specified) with other's index, and sort it. lexicographically.
- inner: form intersection of calling frame's index (or column if on is specified) with other's index, preserving the order of the calling's one.

lsuffix: str, default "

Suffix to use from left frame's overlapping columns.

rsuffix : str, default "

Suffix to use from right frame's overlapping columns.

sort : bool, default False

Order result DataFrame lexicographically by the join key. If False, the order of the join key depends on the join type (how keyword).



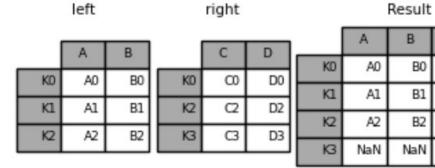
Pandas – join()

result = left.join(right)

result = left.join(right, how="outer")

result = left.join(right, how="inner")

| | left | | | right | | | | Result | | | |
|----|------|----|----|-------|----|----|----|--------|-----|-----|--|
| | Α | В | | С | D | | Α | В | С | D | |
| K0 | A0 | В0 | KO | œ | D0 | KO | A0 | B0 | œ | D0 | |
| K1 | Al | B1 | K2 | C2 | D2 | K1 | Al | B1 | NaN | NaN | |
| K2 | A2 | B2 | КЗ | СЗ | D3 | K2 | A2 | B2 | C2 | D2 | |

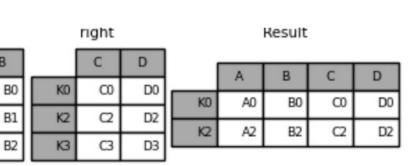


left

A2

K1

K2



B0

NaN

 $^{\circ}$

NaN

C2

C3

D0

NaN

D2

D3



Pandas – join() – example: left / right

Join DataFrames





Pandas – join() – example: inner / outer

```
In [16]:
          df join inner = df1.join(df3, how = 'inner', lsuffix='df1')
          df join inner
Out[16]:
                        namedf1
                                       date height weight age
                                                                                     status children sex
                                                                       name hours
                 id
           1373913 Marisa Martins 2013-02-05
                                               155
                                                            45 Marisa Martins
                                                                                 3 married
                                                                                                  2
                                                       48
          df_join_outer = df1.join(df3, how = 'outer', lsuffix='df1')
          df join outer
Out[17]:
                        namedf1
                                       date height weight age
                                                                        name hours
                                                                                      status children
                id
           1109818
                     Rita Fonseca 2018-08-28
                                             166.0
                                                      54.0
                                                           45.0
                                                                                                NaN NaN
                                                                         NaN
                                                                                NaN
                                                                                        NaN
           1158813
                                       NaT
                                                                  Joana Freitas
                                                                                                  1.0
                                                                                                         F
                            NaN
                                              NaN
                                                      NaN
                                                           NaN
                                                                                 3.0
                                                                                      widow
           1373913
                    Marisa Martins 2013-02-05
                                             155.0
                                                      48.0
                                                           45.0
                                                                 Marisa Martins
                                                                                                  2.0
                                                                                                         F
                                                                                 3.0
                                                                                     married
           1767703
                                                                Manuel Martins
                            NaN
                                       NaT
                                              NaN
                                                      NaN NaN
                                                                                 6.0
                                                                                       single
                                                                                                  0.0
                                                                                                        M
```