# Merge, join, and concatenate

pandas provides various facilities for easily combining together Series, DataFrame, and Panel objects with various kinds of set logic for the indexes and relational algebra functionality in the case of join / merge-type operations.

## **Concatenating objects**

The concat function (in the main pandas namespace) does all of the heavy lifting of performing concatenation operations along an axis while performing optional set logic (union or intersection) of the indexes (if any) on the other axes. Note that I say "if any" because there is only a single possible axis of concatenation for Series.

Before diving into all of the details of concat and what it can do, here is a simple example:

```
In [1]: df1 = pd. DataFrame({'A': ['AO', 'A1', 'A2', 'A3'],
                                                   , 'B2', bu . 'C3'],
                                'B': ['B0', 'B1', 'B2', 'B3'],
'C': ['C0', 'C1', 'C2', 'C3'],
'D': ['D0', 'D1', 'D2', 'D3']},
   . . . :
                                index=[0, 1, 2, 3])
   . . . :
In [2]: df2 = pd. DataFrame({'A': ['A4', 'A5', 'A6', 'A7'],
                                'B': ['B4', 'B5', 'B6', 'B7'],
'C': ['C4', 'C5', 'C6', 'C7'],
   . . . :
                                'D': ['D4', 'D5', 'D6', 'D7']},
                                 index=[4, 5, 6, 7])
   . . . :
   . . . :
'D': ['D8', 'D9', 'D10', 'D11']},
                                index=[8, 9, 10, 11])
   . . . :
In [4]: frames = [df1, df2, df3]
In [5]: result = pd. concat(frames)
```

		df1					Result		
	Α	В	С	D					
0	A0	В0	co	D0		Α	В	С	D
1	A1	B1	C1	D1	0	A0	В0	8	D0
2	A2	B2	C2	D2	1	A1	B1	C1	D1
3	A3	В3	C3	D3	2	A2	B2	C2	D2
		df2							
	Α	В	С	D	3	A3	B3	СЗ	D3
4	A4	B4	C4	D4	4	A4	B4	C4	D4
5	A5	B5	C5	D5	5	A5	B5	C5	D5
6	Аб	В6	C6	D6	6	Aб	В6	C6	D6
7	A7	В7	C7	D7	7	A7	B7	C7	D7
		df3							
	Α	В	С	D	8	A8	B8	C8	DB
8	A8	B8	C8	DB	9	A9	B9	C9	D9
9	A9	B9	C9	D9	10	A10	B10	C10	D10
10	A10	B10	C10	D10	11	A11	B11	C11	D11
11	A11	B11	C11	D11					

Like its sibling function on ndarrays, numpy. concatenate, pandas. concat takes a list or dict of homogeneously-typed objects and concatenates them with some configurable handling of "what to do with the other axes":

- objs: list or dict of Series, DataFrame, or Panel objects. If a dict is passed, the sorted keys will be used as the *keys* argument, unless it is passed, in which case the values will be selected (see below)
- axis: {0, 1, ...}, default 0. The axis to concatenate along
- join: {'inner', 'outer'}, default 'outer'. How to handle indexes on other axis(es). Outer for union and inner for intersection
- join\_axes: list of Index objects. Specific indexes to use for the other n 1 axes instead of performing inner/outer set logic
- keys: sequence, default None. Construct hierarchical index using the passed keys as the outermost level If multiple levels passed, should contain tuples.
- levels: list of sequences, default None. If keys passed, specific levels to use for the resulting MultiIndex. Otherwise they will be inferred from the keys
- names: list, default None. Names for the levels in the resulting hierarchical index
- verify\_integrity: boolean, default False. Check whether the new concatenated axis contains duplicates. This can be very expensive relative to the actual data concatenation
- ignore\_index : boolean, default False. If True, do not use the index values on the concatenation axis. The resulting axis will be labeled 0, ..., n 1. This is useful if you are concatenating objects where the concatenation axis does not have meaningful indexing information.

Without a little bit of context and example many of these arguments don't make much sense. Let's take the above example. Suppose we wanted to associate specific keys with each of the pieces of the chopped up DataFrame. We can do this using the keys argument:

```
In [6]: result = pd. concat(frames, keys=['x', 'y', 'z'])
```

		df1					Res	sult		
	Α	В	С	D						
0	A0	В0	co	D0			Α	В	С	D
1	Al	B1	C1	D1	х	0	A0	B0	α	D0
2	A2	B2	C2	D2	×	1	A1	B1	CI.	D1
3	A3	В3	СЗ	D3	×	2	A2	B2	02	D2
		df2							_	
	A B C D					3	A3	B3	СЗ	D3
4	A4	B4	C4	D4	у	4	A4	B4	C4	D4
5	A5	B5	C5	D5	у	5	A5	B5	C5	D5
6	A6	B6	C6	D6	у	6	Аб	B6	O6	D6
7	A7	В7	C7	D7	у	7	A7	B7	C7	D7
		df3								
	Α	В	С	D	z	8	A8	B8	C8	D8
8	A8	B8	C8	DB	z	9	A9	B9	C9	D9
9	A9	B9	C9	D9	z	10	A10	B10	C10	D10
10	A10	B10	C10	D10	z	11	A11	B11	CI1	D11
11	A11	B11	C11	D11						

As you can see (if you've read the rest of the documentation), the resulting object's index has a *hierarchical index*. This means that we can now do stuff like select out each chunk by key:

```
In [7]: result.ix['y']
Out[7]:
    A   B   C   D
4   A4   B4   C4   D4
5   A5   B5   C5   D5
6   A6   B6   C6   D6
7   A7   B7   C7   D7
```

It's not a stretch to see how this can be very useful. More detail on this functionality below.

**Note:** It is worth noting however, that concat (and therefore append) makes a full copy of the data, and that constantly reusing this function can create a significant performance hit. If you need to use the operation over several datasets, use a list comprehension.

```
frames = [ process_your_file(f) for f in files ]
result = pd. concat(frames)
```

### Set logic on the other axes

When gluing together multiple DataFrames (or Panels or...), for example, you have a choice of

how to handle the other axes (other than the one being concatenated). This can be done in three ways:

- Take the (sorted) union of them all, join='outer'. This is the default option as it results in zero information loss.
- Take the intersection, join='inner'.
- Use a specific index (in the case of DataFrame) or indexes (in the case of Panel or future higher dimensional objects), i.e. the join\_axes argument

Here is a example of each of these methods. First, the default <code>join='outer'</code> behavior:

		df1				df	4					Res	ult			
										Α	В	С	D	В	D	F
	Α	В	С	D		В	D	F	0	A0	В0	CO	D0	NaN	NaN	NaN
0	A0	В0	œ	D0	2	B2	D2	F2	1	A1	В1	C1	D1	NaN	NaN	NaN
1	A1	B1	C1	D1	3	В3	D3	F3	2	A2	B2	C2	D2	B2	D2	F2
2	A2	B2	C2	D2	6	B6	D6	F6	3	A3	В3	C3	D3	В3	D3	F3
3	A3	В3	C3	D3	7	В7	D7	F7	6	NaN	NaN	NaN	NaN	В6	D6	F6
									7	NaN	NaN	NaN	NaN	В7	D7	F7

Note that the row indexes have been unioned and sorted. Here is the same thing with join=' inner':

```
In [10]: result = pd.concat([df1, df4], axis=1, join='inner')
```

		df1				df	4					Res	ult			
	Α	В	С	D		В	D	F								
0	A0	BO	co	D0	2	B2	D2	F2		Α	В	С	D	В	D	F
1	A1	B1	Cl	D1	3	В3	D3	F3	2	A2	B2	C2	D2	B2	D2	F2
2	A2	B2	C2	D2	6	В6	D6	F6	3	A3	В3	C3	D3	В3	D3	F3
3	A3	В3	C3	D3	7	В7	D7	F7								

Lastly, suppose we just wanted to reuse the *exact index* from the original DataFrame:

In [11]: result = pd.concat([df1, df4], axis=1, join\_axes=[df1.index])

			df1				df	F4					Res	ult			
		Α	В	С	D		В	D	F		Α	В	С	D	В	D	F
	0	A0	BO	O	D0	2	B2	D2	F2	0	A0	В0	CO	D0	NaN	NaN	NaN
	1	A1	B1	Cl	D1	3	В3	D3	F3	1	A1	В1	C1	D1	NaN	NaN	NaN
Ī	2	A2	B2	(2	D2	6	В6	D6	F6	2	A2	B2	C2	D2	B2	D2	F2
Ī	3	A3	В3	СЗ	D3	7	B7	D7	F7	3	A3	В3	C3	D3	В3	D3	F3

## Concatenating using append

A useful shortcut to concat are the append instance methods on Series and DataFrame. These methods actually predated concat. They concatenate along axis=0, namely the index:

In [12]: result = df1.append(df2)

		df1					Result		
	Α	В	С	D		Α	В	С	D
0	A0	В0	00	D0			DO.		- Do
1	A1	B1	C1	D1	0	A0	B0	00	D0
2	A2	B2	C2	D2	1	A1	B1	C1	D1
3	A3	В3	C3	D3	2	A2	B2	C2	D2
		df2			3	A3	В3	СЗ	D3
	Α	В	С	D	4	A4	B4	C4	D4
4	A4	B4	C4	D4	5	A5	B5	C5	D5
5	A5	B5	C5	D5		45	D.C.	~	D.C.
6	A6	В6	O6	D6	6	A6	B6	O6	D6
7	A7	B7	C7	D7	7	A7	В7	C7	D7

In the case of DataFrame, the indexes must be disjoint but the columns do not need to be:

```
In [13]: result = df1.append(df4)
```

			df	F1					Res	sult		
		Α	E	3	С	D		Α	В	С	D	F
C	1	A0		BO	α	D0	_	40	DO.		D.O.	NI - NI
1		A1		В1	C	1 D1	0	A0	В0	co	D0	NaN
- 2	2	A2		В2	C	2 D2	1	A1	В1	C1	D1	NaN
3	3	А3		ВЗ	C	3 D3	2	A2	В2	C2	D2	NaN
		df4				3	А3	В3	C3	D3	NaN	
		В		ı	D	F	2	NaN	В2	NaN	D2	F2
	2	ı	В2		D2	F2	3	NaN	В3	NaN	D3	F3
	3		ВЗ		D3	F3	6	NaN	D.C	NISNI	De	F6
	6		В6		D6	F6	6	NaN	В6	NaN	D6	F6
	7		В7		D7	F7	7	NaN	В7	NaN	D7	F7

append may take multiple objects to concatenate:

		df1					Result		
	Α	В	С	D					
0	A0	BO	œ	D0		Α	В	С	D
1	A1	B1	C1	D1	0	A0	В0	α	D0
2	A2	B2	C2	D2	1	A1	B1	C1	D1
3	A3	В3	C3	D3	2	A2	B2	C2	D2
		df2			_				
	Α	В	С	D	3	A3	В3	C3	D3
4	A4	B4	C4	D4	4	A4	B4	C4	D4
5	A5	B5	C5	D5	5	A5	B5	C5	D5
6	A6	B6	C6	D6	6	A6	В6	C6	D6
7	A7	B7	C7	D7	7	A7	В7	C7	D7
		df3							
	Α	В	С	D	8	A8	B8	C8	DB
8	A8	B8	C8	DB	9	A9	B9	C9	D9
9	A9	B9	C9	D9	10	A10	B10	C10	D10
10	A10	B10	C10	D10	11	A11	B11	C11	D11
11	A11	B11	C11	D11					

**Note:** Unlike *list.append* method, which appends to the original list and returns nothing, append here **does not** modify df1 and returns its copy with df2 appended.

## Ignoring indexes on the concatenation axis

For DataFrames which don't have a meaningful index, you may wish to append them and ignore the fact that they may have overlapping indexes:

To do this, use the <code>ignore\_index</code> argument:

```
In [15]: result = pd.concat([df1, df4], ignore_index=True)
```

			df	1					Res	sult		
		Α	E	3	С	D		Α	В	С	D	F
	0	A0		BO	0	0 D0	_					
Г	1	A1		B1	С	1 D1	0	A0	В0	CO	D0	NaN
	2	A2		B2	C	2 D2	1	A1	В1	C1	D1	NaN
	3	A3		ВЗ	C	3 D3	2	A2	В2	C2	D2	NaN
Ī			df	4			3	А3	В3	C3	D3	NaN
_		В			D	F	4	NaN	В2	NaN	D2	F2
	- 2	2	В2		D2	F2	5	NaN	В3	NaN	D3	F3
		3	ВЗ		D3	F3	6	NaN	В6	NaN	D6	F6
	(	5	В6		D6	F6	0	INGIN	00	INGIN	D6	го
		7	В7		D7	F7	7	NaN	В7	NaN	D7	F7

This is also a valid argument to DataFrame. append:

		df1						Res	sult		
	Α	В	С	D			Α	В	С	D	F
0	A0	В	0 0	0D D0	lı	_	40	DO.			N1 - N1
1	A1	В	1 C	1 D1	11	0	A0	В0	CO	D0	NaN
2	A2	В	2 0	2 D2	1	1	A1	В1	C1	D1	NaN
3	A3	В	3 C	3 D3	1	2	A2	В2	C2	D2	NaN
		df4			1	3	А3	В3	C3	D3	NaN
	В		D	F		4	NaN	В2	NaN	D2	F2
- 2	2	B2	D2	F2		5	NaN	В3	NaN	D3	F3
3	3	В3	D3	F3	Н	6	NaN	В6	NaN	D6	F6
(	5	B6	D6	F6	11		Ivalv	Во	Ivalv		
	7	В7	D7	F7	$\ $	7	NaN	В7	NaN	D7	F7

## Concatenating with mixed ndims

You can concatenate a mix of Series and DataFrames. The Series will be transformed to DataFrames with the column name as the name of the Series.

```
In [17]: s1 = pd. Series(['X0', 'X1', 'X2', 'X3'], name='X')
In [18]: result = pd. concat([df1, s1], axis=1)
```

		df1			S	1			Res	ult		
	Α	В	С	D		Х		Α	В	С	D	Х
0	A0	В0	œ	D0	0	X0	0	A0	В0	CO	D0	X0
1	A1	B1	Cl	D1	1	X1	1	A1	В1	C1	D1	X1
2	A2	B2	C2	D2	2	X2	2	A2	B2	C2	D2	X2
3	A3	В3	C3	D3	3	ХЗ	3	A3	В3	C3	D3	ХЗ

If unnamed Series are passed they will be numbered consecutively.

```
In [19]: s2 = pd. Series(['_0', '_1', '_2', '_3'])
In [20]: result = pd. concat([df1, s2, s2, s2], axis=1)
```

		df1			S	2				Res	sult			
	Α	В	С	D				Α	В	С	D	0	1	2
0	A0	BO	α	D0	0	_0	0	A0	В0	co	D0	_0	_0	_0
1	A1	B1	C1	D1	1	_1	1	A1	В1	C1	D1	_1	_1	_1
2	A2	B2	C2	D2	2	_2	2	A2	В2	C2	D2	_2	_2	_2
3	A3	В3	C3	D3	3	_3	3	A3	В3	СЗ	D3	_3	_3	_3

Passing ignore\_index=True will drop all name references.

```
In [21]: result = pd. concat([df1, s1], axis=1, ignore_index=True)
```

		df1			S	1			Res	ult		
	Α	В	С	D		Х		0	1	2	3	4
0	A0	В0	α	D0	0	X0	0	A0	В0	CO	D0	X0
1	A1	B1	CI	D1	1	X1	1	A1	В1	C1	D1	X1
2	A2	B2	C2	D2	2	X2	2	A2	В2	C2	D2	X2
3	A3	В3	C3	D3	3	ХЗ	3	A3	В3	C3	D3	ХЗ

#### More concatenating with group keys

A fairly common use of the keys argument is to override the column names when creating a new DataFrame based on existing Series. Notice how the default behaviour consists on letting the resulting DataFrame inherits the parent Series' name, when these existed.

```
In [22]: s3 = pd. Series([0, 1, 2, 3], name='foo')
In [23]: s4 = pd. Series([0, 1, 2, 3])
In [24]: s5 = pd. Series([0, 1, 4, 5])
In [25]: pd. concat([s3, s4, s5], axis=1)
Out[25]:
    foo 0 1
0 0 0 0
1 1 1 1 1
2 2 2 2 4
3 3 3 5
```

Through the keys argument we can override the existing column names.

```
In [26]: pd.concat([s3, s4, s5], axis=1, keys=['red', 'blue', 'yellow'])
Out [26]:
  red blue yellow
0
     0
           0
     1
           1
1
                    1
2
     2
           2
                    4
     3
           3
3
                    5
```

Let's consider now a variation on the very first example presented:

```
In [27]: result = pd. concat(frames, keys=['x', 'y', 'z'])
```

		df1					Res	sult		
	Α	В	С	D						
0	A0	В0	co	D0			А	В	С	D
1	Al	B1	C1	D1	х	0	A0	B0	α	D0
2	A2	B2	C2	D2	×	1	Al	B1	a	D1
3	A3	В3	C3	D3	×	2	A2	B2	(2	D2
		df2								
	Α	В	С	D	х	3	A3	B3	G	D3
4	A4	B4	C4	D4	у	4	A4	B4	C4	D4
5	A5	B5	C5	D5	у	5	A5	B5	CS	D5
6	Аб	B6	C6	D6	у	6	Аб	B6	C6	D6
7	A7	B7	C7	D7	у	7	A7	B7	C7	D7
		df3								
	Α	В	С	D	Z	8	A8	B8	C8	D8
8	A8	B8	C8	DB	z	9	A9	B9	C9	D9
9	A9	B9	C9	D9	z	10	A10	B10	C10	D10
10	A10	B10	C10	D10	z	11	A11	B11	CI1	D11
11	11 All Bll Cll Dll									

You can also pass a dict to  ${\tt concat}$  in which case the dict keys will be used for the  ${\tt keys}$  argument (unless other keys are specified):

```
In [28]: pieces = {'x': df1, 'y': df2, 'z': df3}
In [29]: result = pd. concat(pieces)
```

		df1					Res	sult		
	Α	В	С	D						
0	A0	В0	co	D0			Α	В	С	D
1	Al	B1	C1	D1	х	0	A0	B0	Θ	D0
2	A2	B2	C2	D2	×	1	Al	B1	a	D1
3	A3	В3	C3	D3	×	2	A2	B2	(2	D2
		df2								
	Α	В	С	D	х	3	A3	B3	G	D3
4	A4	B4	C4	D4	у	4	A4	B4	C4	D4
5	A5	B5	C5	D5	у	5	A5	B5	CS	D5
6	Аб	B6	O6	D6	у	6	Аб	B6	C6	D6
7	A7	B7	C7	D7	у	7	A7	B7	C7	D7
		df3								
	Α	В	С	D	Z	8	A8	B8	C8	D8
8	A8	B8	C8	DB	Z	9	A9	B9	(9	D9
9	A9	B9	C9	D9	z	10	A10	B10	C10	D10
10	A10	B10	C10	D10	z	11	A11	B11	Cll	D11
11	A11	B11	C11	D11						

```
In [30]: result = pd. concat(pieces, keys=['z', 'y'])
```

								_			
			df1					Res	sult		
		Α	В	С	D						
	0	A0	В0	co	D0						
	1	Al	B1	C1	D1						
	2	A2	B2	C2	D2			А	В	С	D
	3	A3	В3	C3	D3	z	8	AB	B8	08	D8
			df2								
		Α	В	С	D	Z	9	A9	B9	(9	D9
	4	A4	B4	C4	D4	z	10	A10	B10	C10	D10
ı	5	A5	B5	C5	D5	z	11	A11	B11	C11	D11
	6	A6	B6	C6	D6	у	4	A4	B4	C4	D4
	7	A7	B7	C7	D7	у	5	A5	B5	C5	D5
			df3								$\vdash$
		Α	В	С	D	У	6	Аб	B6	06	D6
	8	A8	B8	C8	DB	У	7	A7	B7	C7	D7
	9	A9	B9	C9	D9						
ı	10	A10	B10	C10	D10						
	11	A11	B11	C11	D11						

The MultiIndex created has levels that are constructed from the passed keys and the index of the DataFrame pieces:

```
In [31]: result.index.levels
Out[31]: FrozenList([[u'z', u'y'], [4, 5, 6, 7, 8, 9, 10, 11]])
```

If you wish to specify other levels (as will occasionally be the case), you can do so using the levels argument:

		df1					Res	sult		
	Α	В	С	D						
0	A0	В0	CO	D0			Α	В	С	D
1	Al	B1	C1	D1	х	0	A0	В0	CO	D0
2	A2	B2	C2	D2	х	1	A1	B1	CI	D1
3	A3	В3	C3	D3	×	2	A2	B2	(2	D2
		df2								
	Α	В	С	D	х	3	A3	B3	G	D3
4	A4	B4	C4	D4	У	4	A4	B4	C4	D4
5	A5	B5	C5	D5	у	5	A5	B5	C5	D5
6	A6	B6	C6	D6	у	6	Аб	B6	C6	D6
7	A7	B7	C7	D7	у	7	A7	B7	C7	D7
		df3					40			
	Α	В	С	D	Z	8	A8	B8	C8	D8
8	A8	B8	C8	DB	z	9	A9	B9	C9	D9
9	A9	B9	C9	D9	Z	10	A10	B10	C10	D10
10	A10	B10	C10		z	11	A11	B11	Cll	D11
11	A11	All Bll Cll Dll	D11							

```
In [33]: result.index.levels
Out[33]: FrozenList([[u'z', u'y', u'x', u'w'], [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]])
```

Yes, this is fairly esoteric, but is actually necessary for implementing things like GroupBy where the order of a categorical variable is meaningful.

## Appending rows to a DataFrame

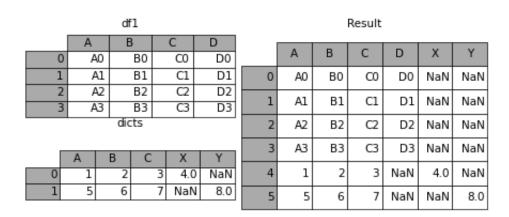
While not especially efficient (since a new object must be created), you can append a single row to a DataFrame by passing a Series or dict to append, which returns a new DataFrame as above.

```
In [34]: s2 = pd. Series(['X0', 'X1', 'X2', 'X3'], index=['A', 'B', 'C', 'D'])
In [35]: result = df1.append(s2, ignore_index=True)
```

			df1					Result		
		Α	В	С	D					
	0	A0	В0	α	D0					
	1	Al	B1	C1	D1		Α	В	С	D
l	2	A2	B2	(2	D2	0	A0	BO	œ	D0
l	3	A3	В3	СЗ	D3	1	Al	B1	Cl	D1
			s2							
						2	A2	B2	C2	D2
			Α		XO	3	A3	В3	СЗ	D3
İ			В		X1	4	X0	X1	X2	ХЗ
Ī			С		X2					
Ī			D		ХЗ					

You should use <code>ignore\_index</code> with this method to instruct DataFrame to discard its index. If you wish to preserve the index, you should construct an appropriately-indexed DataFrame and append or concatenate those objects.

You can also pass a list of dicts or Series:



## Database-style DataFrame joining/merging

pandas has full-featured, **high performance** in-memory join operations idiomatically very similar to relational databases like SQL. These methods perform significantly better (in some cases well over an order of magnitude better) than other open source implementations (like base::merge. data. frame in R). The reason for this is careful algorithmic design and internal layout of the data in DataFrame.

See the *cookbook* for some advanced strategies.

Users who are familiar with SQL but new to pandas might be interested in a *comparison with SQL*.

pandas provides a single function, merge, as the entry point for all standard database join operations between DataFrame objects:

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

Here's a description of what each argument is for:

- left: A DataFrame object
- right: Another DataFrame object
- on: Columns (names) to join on. Must be found in both the left and right
   DataFrame objects. If not passed and left\_index and right\_index are False,
   the intersection of the columns in the DataFrames will be inferred to be the
   join keys
- left\_on: Columns from the left DataFrame to use as keys. Can either be column names or arrays with length equal to the length of the DataFrame
- right\_on: Columns from the right DataFrame to use as keys. Can either be column names or arrays with length equal to the length of the DataFrame
- left\_index: If True, use the index (row labels) from the left DataFrame as its join key(s). In the case of a DataFrame with a MultiIndex (hierarchical), the number of levels must match the number of join keys from the right DataFrame
- right\_index: Same usage as left\_index for the right DataFrame
- how: One of 'left', 'right', 'outer', 'inner'. Defaults to inner. See below for more detailed description of each method
- 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 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, right\_only for observations whose merge key only appears in 'right' DataFrame, and both if the observation's merge key is found in both.

New in version 0.17.0.

The return type will be the same as left. If left is a DataFrame and right is a subclass of DataFrame, the return type will still be DataFrame.

merge is a function in the pandas namespace, and it is also available as a DataFrame instance method, with the calling DataFrame being implicitly considered the left object in the join.

The related DataFrame. join method, uses merge internally for the index-on-index and index-on-column(s) joins, but *joins on indexes* by default rather than trying to join on common columns (the default behavior for merge). If you are joining on index, you may wish to use DataFrame. join to save yourself some typing.

#### Brief primer on merge methods (relational algebra)

Experienced users of relational databases like SQL will be familiar with the terminology used to describe join operations between two SQL-table like structures (DataFrame objects). There are several cases to consider which are very important to understand:

- **one-to-one** joins: for example when joining two DataFrame objects on their indexes (which must contain unique values)
- many-to-one joins: for example when joining an index (unique) to one or more columns in a DataFrame
- many-to-many joins: joining columns on columns.

**Note:** When joining columns on columns (potentially a many-to-many join), any indexes on the passed DataFrame objects **will be discarded**.

It is worth spending some time understanding the result of the **many-to-many** join case. In SQL / standard relational algebra, if a key combination appears more than once in both tables, the resulting table will have the **Cartesian product** of the associated data. Here is a very basic example with one unique key combination:

	le	ft			rig	ht				Res	sult		
A B key					С	D	key		Α	В	key	С	D
0	A0	В0	KO	0	α	D0	K0	0	A0	В0	K0	CO	D0
1	A1	B1	K1	1	CI	D1	K1	1	A1	В1	K1	C1	D1
2	A2	B2	K2	2	C2	D2	K2	2	A2	B2	K2	C2	D2
3	A3	В3	КЗ	3	СЗ	D3	КЗ	3	A3	В3	К3	C3	D3

Here is a more complicated example with multiple join keys:

			left					right						Result			
		Α	В	keyl	key2		С	D	keyl	key2		Α	В	key1	kay2	С	D
	0	A0	В0	K0	K0	0	α	D0	KO	K0					key2		
1	1	A1	B1	KO	К1	1	C1	D1	K1	KO	0	A0	B0	K0	K0	co	D0
		~1	DI		$\vdash$						1	A2	B2	K1	K0	C1	D1
	2	A2	B2	K1	K0	2	(2	D2	K1	KO	2	42					Da
ı	3	A3	В3	K2	К1	3	СЗ	D3	K2	KO		A2	B2	K1	K0	C2	D2
- 1																	

The how argument to merge specifies how to determine which keys are to be included in the resulting table. If a key combination **does not appear** in either the left or right tables, the values in the joined table will be NA. Here is a summary of the how options and their 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

```
In [44]: result = pd.merge(left, right, how='left', on=['key1', 'key2'])
```

left right key1 key2 key1  $\alpha$ D0 A2 В2 K1 K0 C2 D2 ΚO АЗ В3 K1 C3 D3 ΚO K2

	Result													
	A B key1 key2 C D													
(	0	A0	В0	K0	K0	CO	D0							
	1	A1	В1	K0	K1	NaN	NaN							
	2	A2	В2	K1	K0	C1	D1							
	3	A2	B2	K1	K0	C2	D2							
4	4	A3	В3	K2	K1	NaN	NaN							

In [45]: result = pd.merge(left, right, how='right', on=['key1', 'key2'])

			left					right						Result			
		Α	В	keyl	key2		С	D	keyl	key2		Α	В	key1	key2	С	D
	0	A0	BO	K0	K0	0	ω	D0	K0	K0	0	A0	В0	K0	K0	co	D0
	1	A1	B1	K0	K1	1	Cl	D1	K1	K0	1	A2	В2	K1	K0	C1	D1
	2	A2	B2	K1	K0	2	(2	D2	K1	K0	2	A2	В2	K1	K0	C2	D2
I	3	A3	В3	K2	K1	3	СЗ	D3	K2	K0	3	NaN	NaN	K2	K0	C3	D3

In [46]: result = pd.merge(left, right, how='outer', on=['key1', 'key2'])

	left					right						Result			
										Α	В	key1	key2	С	D
Α	В	keyl	key2		С	D	keyl	key2	0	A0	В0	K0	K0	co	D0
A0	BO	K0	KO	0	0	D0	K0	KO	1	A1	В1	K0	K1	NaN	NaN
Al	B1	K0	K1	1	Cl	D1	K1	KO	2	A2	В2	K1	K0	C1	D1
A2	B2	K1	KO	2	C2	D2	K1	KO	3	A2	В2	K1	K0	C2	D2
A3	В3	K2	K1	3	C3	D3	K2	KO	4	A3	В3	K2	K1	NaN	NaN
									5	NaN	NaN	K2	K0	C3	D3
	A0 A1 A2	A B A0 B0 A1 B1 A2 B2	A B key1  A0 B0 K0  A1 B1 K0  A2 B2 K1	A B key1 key2  A0 B0 K0 K0  A1 B1 K0 K1  A2 B2 K1 K0	A B key1 key2  A0 B0 K0 K0 0  A1 B1 K0 K1 1  A2 B2 K1 K0 2	A B key1 key2 C  A0 B0 K0 K0 0 C0  A1 B1 K0 K1 1 C1  A2 B2 K1 K0 2 C2	A B keyl key2 C D  A0 B0 K0 K0 0 C0 D0  A1 B1 K0 K1 1 C1 D1  A2 B2 K1 K0 2 C2 D2	A B key1 key2 C D key1  A0 B0 K0 K0 0 C0 D0 K0  A1 B1 K0 K1 1 C1 D1 K1  A2 B2 K1 K0 2 C2 D2 K1	A         B         key1         key2         C         D         key1         key2           A0         B0         K0         K0         0         C0         D0         K0         K0           A1         B1         K0         K1         1         C1         D1         K1         K0           A2         B2         K1         K0         2         C2         D2         K1         K0	A B keyl key2	A B keyl key2 C D keyl key2 0 A0 A0 B0 K0 K0 0 0 00 D0 K0 K0 1 A1 A1 B1 K0 K1 1 C1 D1 K1 K0 2 A2 A2 B2 K1 K0 2 C2 D2 K1 K0 3 A2 A3 B3 K2 K1 3 C3 D3 K2 K0 4 A3	A         B         key1         key2         C         D         key1         key2         0         A0         B0           A0         B0         K0         K0         0         0         D0         K0         K0         1         A1         B1           A1         B1         K0         K1         1         C1         D1         K1         K0         2         A2         B2           A2         B2         K1         K0         2         C2         D2         K1         K0         3         A2         B2           A3         B3         K2         K1         3         C3         D3         K2         K0         4         A3         B3	A B key1 A B key1 key2 C D key1 key2 O A0 B0 K0 A0 B0 K0 K0 O C0 D0 K0 K0 A1 B1 K0 K1 1 C1 D1 K1 K0 2 A2 B2 K1 A3 B3 K2 K1 3 C3 D3 K2 K0 A B key1 A0 B0 key1 A0 B0 K0 A1 B1 K0 A1 B1 K0 A1 B1 K0 A2 B2 K1 A3 B3 K2 K1 A3 B3 K2 K1 A3 B3 K2 K1	A B key1 key2  A B key1 key2  C D key1 key2  O A0 B0 K0 K0 K0  A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1  A3 C3 D3 K2 K0  A B key1 key2  O A0 B0 K0 K0 K0  A1 A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1  A B key1 key2  O A0 B0 K0 K0  A0 B0 K0 K0  A1 A1 B1 K0 K1  A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1  A8 B key1 key2  O A0 B0 K0 K0  A0 K1  A1 B1 K0 K1  A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1  A8 B8 key1 key2  O A0 B0 K0 K0  A0 K1  A1 B1 K0 K1  A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1  A8 B8 key1 key2  O A0 B0 K0 K0  A0 K1  A1 B1 K0 K1  A1 B1 K0 K1  A1 B1 K0 K1  A2 B2 K1 K0  A3 B3 K2 K1	A B key1 key2 C  A B key1 key2 C  C D key1 key2 0 A0 B0 K0 K0 K0 C0  A0 B0 K0 K0 0 0 C0 D0 K0 K0 1 A1 B1 K0 K1 NaN  A1 B1 K0 K1 1 C1 D1 K1 K0 2 A2 B2 K1 K0 C1  A2 B2 K1 K0 2 C2 D2 K1 K0 3 A2 B2 K1 K0 C2  A3 B3 K2 K1 3 C3 D3 K2 K0 4 A3 B3 K2 K1 NaN

```
In [47]: result = pd.merge(left, right, how='inner', on=['key1', 'key2'])
```

	left						right			Result						
	Α	В	keyl	key2		С	D	keyl	key2		Α	В	key1	kay2	С	D
0	A0	В0	K0	K0	0	α	D0	K0	K0							
1	A1	B1	KO	К1	1	CI	D1	K1	KO	0	A0	В0	K0	K0	co	D0
						-				1	A2	B2	K1	K0	C1	D1
2	A2	B2	K1	K0	2	(2	D2	K1	K0	2	A2	B2	K1	К0	C2	D2
3	A3	В3	K2	K1	3	C3	D3	K2	KO		72	02	14.1	100	- 02	02
_									1.00							

#### The merge indicator

New in version 0.17.0.

merge now accepts the argument indicator. If True, a Categorical-type column called \_merge will be added to the output object that takes on values:

Observation Origin	_merge value
Merge key only in 'left' frame	left_only
Merge key only in 'right' frame	right_only
Merge key in both frames	both

```
In [48]: df1 = DataFrame({'col1':[0,1], 'col left':['a','b']})
In [49]: df2 = DataFrame({'coll':[1,2,2],'col right':[2,2,2]})
In [50]: merge(df1, df2, on='col1', how='outer', indicator=True)
Out [50]:
   coll col left col right
                                 merge
0
                        NaN left only
               а
                          2
1
      1
               b
                                   both
2
      2
                          2 right_only
             NaN
3
      2
             NaN
                          2 right_only
```

The indicator argument will also accept string arguments, in which case the indicator function will use the value of the passed string as the name for the indicator column.

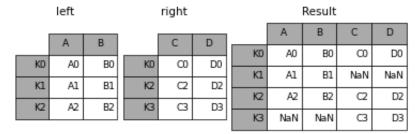
```
In [51]: merge(df1, df2, on='coll', how='outer', indicator='indicator_column')
Out[51]:
   coll col_left col_right indicator_column
0
      0
                        NaN
                                    left only
               а
1
      1
               b
                           2
                                         both
2
      2
                           2
                                   right only
             NaN
3
             NaN
                                   right_only
```

### Joining on index

DataFrame. join is a convenient method for combining the columns of two potentially differently-indexed DataFrames into a single result DataFrame. Here is a very basic example:

	left			right				Result		
	Α	В		С	D		Α	В	С	D
K0	A0	В0	KO	8	D0	KO	A0	В0	8	D0
Κ1	A1	B1	K2	C2	D2	K1	A1	B1	NaN	NaN
K2	A2	B2	КЗ	C3	D3	K2	A2	B2	C2	D2

```
In [55]: result = left.join(right, how='outer')
```



```
In [56]: result = left.join(right, how='inner')
```

	left			right			
	Α	В		С	D		
KO	A0	BO	KO	α	D0	KO	
K1	A1	B1	K2	C2	D2	K2	
K2	A2	B2	КЗ	СЗ	D3	NZ.	

	Α	В	С	D
K0	A0	В0	œ	D0
K2	A2	B2	C2	D2

Result

The data alignment here is on the indexes (row labels). This same behavior can be achieved using merge plus additional arguments instructing it to use the indexes:

```
In [57]: result = pd.merge(left, right, left_index=True, right_index=True, how='outer')
```

		left			right Result						
		Α	В		С	D		Α	В	С	D
							KO	A0	В0	α	D0
ш	K0	A0	B0	KO	00	D0					
н	1/2	43		160			K1	A1	B1	NaN	NaN
ш	K1	A1	B1	K2	C2	D2	K2	A2	B2	C2	D2
	K2	A2	B2	КЗ	C3	D3	1/4	PZ	52	- 02	D2
								NaN	NaN	C3	D3

```
In [58]: result = pd.merge(left, right, left_index=True, right_index=True, how='inner');
```

	left				right		Result					
	Α		В		С	D		Δ	В	С	D	
K	0 A	7	BO	KO	00	D0		^		_	-	
K	1 A	+	B1	K2	C2	D2	KO	A0	B0	α	D0	
,	^	1	DI	I/Z	- 02	102	K2	A2	B2	C2	D2	
K	2 A3	2	B2	КЗ	C3	D3						

## Joining key columns on an index

join takes an optional on argument which may be a column or multiple column names, which specifies that the passed DataFrame is to be aligned on that column in the DataFrame. These two function calls are completely equivalent:

```
left. join(right, on=key_or_keys)
pd. merge(left, right, left_on=key_or_keys, right_index=True,
    how='left', sort=False)
```

Obviously you can choose whichever form you find more convenient. For many-to-one joins (where one of the DataFrame's is already indexed by the join key), using join may be more convenient. Here is a simple example:

```
left
                           right
                                                        Result
    Α
                                                                    С
          В
               key
                                                                          D
                                                             key
                              С
    A0
                                                 A0
                                                        B0
                                                              K0
                                                                     \infty
                                                                           D0
1
                 K1
                                             1
                                                        В1
           В1
                              \alpha
                                    D0
                                                  A1
                                                              K1
                                                                           D1
2
     A2
           B2
                 K0
                        K1
                              C1
                                    D1
                                             2
                                                        B2
                                                              K0
                                                                           D0
                                                  A2
                                                                     ^{\circ}
3
                 K1
     АЗ
                                                  ΑЗ
                                                        ВЗ
                                                              K1
                                                                     C1
                                                                           D1
```

	le	ft		right								
	Α	В	key					Α	В	key	С	D
0	A0	В0	KO		С	D	0	A0	В0	K0	CO	D0
1	Al	B1	K1	KO	α	D0	1	A1	В1	K1	C1	D1
2	A2	B2	KO	K1	Cl	D1	2	A2	B2	K0	CO	D0
3	A3	В3	K1				3	A3	В3	K1	C1	D1

To join on multiple keys, the passed DataFrame must have a MultiIndex:

Now this can be joined by passing the two key column names:

```
In [66]: result = left.join(right, on=['key1', 'key2'])
```

left						right				Result					
	Α	В	keyl	key2			С	D		Α	В	key1	key2	С	D
0	A0	В0	K0	KO	K0	KO	α	D0	0	A0	В0	K0	K0	co	D0
1	A1	B1	K0	K1	K1	KO	а	D1	1	A1	В1	K0	K1	NaN	NaN
2	A2	B2	K1	K0	K2	KO	(2	D2	2	A2	В2	K1	K0	C1	D1
3	A3	В3	K2	K1	K2	КІ	З	D3	3	А3	В3	K2	K1	C3	D3

The default for <code>DataFrame.join</code> is to perform a left join (essentially a "VLOOKUP" operation, for Excel users), which uses only the keys found in the calling <code>DataFrame</code>. Other join types, for example inner join, can be just as easily performed:

```
In [67]: result = left.join(right, on=['key1', 'key2'], how='inner')
```

	left					right					Result				
	Α	В	keyl	key2			С	D		Α	В	key1	key2	С	D
0	A0	В0	K0	K0	K0	K0	α	D0				_			
1	A1	B1	KO	К1	KI	KO	а	D1	0	A0	B0	K0	K0	co	D0
_									2	A2	В2	K1	K0	C1	D1
	A2	B2	K1	K0	K2	K0	(2	D2	3	A3	В3	K2	K1	C3	D3
3	A3	В3	K2	K1	K2	K1	G	D3		,,,		142	142		

As you can see, this drops any rows where there was no match.

### Joining a single Index to a Multi-index

New in version 0.14.0.

You can join a singly-indexed DataFrame with a level of a multi-indexed DataFrame. The level will match on the name of the index of the singly-indexed frame against a level name of the multi-indexed frame.

	left			rig	ht		Result					
	А	В			С	D			А	В	С	D
KO	A0	BO	KO	YO	α	D0	KO	YO	A0	B0	α	D0
K1	A1	B1	KI	Y1	а	D1	KI	Y1	Al	B1	а	D1
K2	A2	B2	K2	Y2	O	D2	K2	Y2	A2	B2	(2	D2
NZ.	72	DZ.	K2	Y3	C	D3	K2	Y3	A2	B2	СЗ	D3

This is equivalent but less verbose and more memory efficient / faster than this.

	left			rig	ht				Result			
	۸	В			С	D			А	В	С	D
V0	Α	BO	K0	YO	8	D0	KO	YO	A0	B0	α	D0
K0	A0		Kl	Y1	а	D1	Kl	Y1	Al	B1	a	D1
K1	A1	B1	K2	Y2	(2	D2	K2	Y2	A2	B2	(2	D2
K2	A2	B2	1/2	Y3	СЗ	D3	K2	Y3	A2	B2	СЗ	D3

## Joining with two multi-indexes

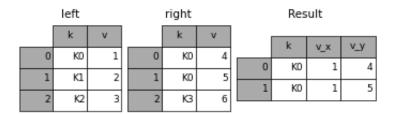
This is not Implemented via join at-the-moment, however it can be done using the following.

left					right				Result						
A B							С	D		A B (					
					KO	70	α	D0						,	D
	K0	XO	A0	BO	100				KD	XO	YO	A0	B0	0	D0
					K1	Y1	a	D1							
	K0	XI	A1	B1	1.2				KD	X1	YO	A1	B1	8	D0
		,4			K2	Y2	0	D2							
	K1	Х2	A2	B2	100		_		K1	Х2	Y1	A2	B2	a	D1
	144	7.2	72		K2	Y3	СЗ	D3	I/C	7/2	12	/2	LL		
							-								

### Overlapping value columns

The merge suffixes argument takes a tuple of list of strings to append to overlapping column names in the input DataFrames to disambiguate the result columns:

```
In [76]: left = pd. DataFrame({'k': ['K0', 'K1', 'K2'], 'v': [1, 2, 3]})
In [77]: right = pd. DataFrame({'k': ['K0', 'K0', 'K3'], 'v': [4, 5, 6]})
In [78]: result = pd. merge(left, right, on='k')
```

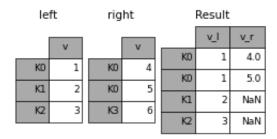


```
In [79]: result = pd. merge(left, right, on='k', suffixes=['_1', '_r'])
```

	left		right				Result				
	k	v		k	v			l l	v I		
0	K0	1	0	K0	4	1		Α.	V_1	v_r	
1	K1	2	1	KO	5	ł	0	KO	1	4	
	11.2			140		ļ	1	KO	1	5	
2	K2	3	2	КЗ	6						

DataFrame. join has 1suffix and rsuffix arguments which behave similarly.

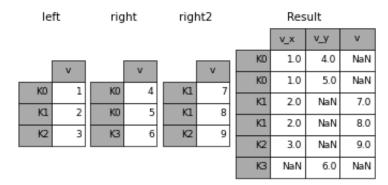
```
In [80]: left = left.set_index('k')
In [81]: right = right.set_index('k')
In [82]: result = left.join(right, lsuffix='_l', rsuffix='_r')
```



## Joining multiple DataFrame or Panel objects

A list or tuple of DataFrames can also be passed to DataFrame. join to join them together on their indexes. The same is true for Panel. join.

```
In [83]: right2 = pd. DataFrame({'v': [7, 8, 9]}, index=['K1', 'K1', 'K2'])
In [84]: result = left.join([right, right2])
```



#### **Merging Ordered Data**

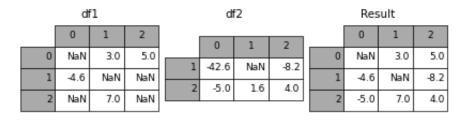
New in v0.8.0 is the ordered\_merge function for combining time series and other ordered data. In particular it has an optional fill\_method keyword to fill/interpolate missing data:

	le	ft			Result							
	k	lv	5	]	k	lv	s	rv				
				0	K0	1.0	a	NaN				
0	KD		l a	1	K1	1.0	a	1.0				
1	K1	:	2 Ь	2	K2	1.0	a	2.0				
2	K1		3 с	3	K4	1.0	a	3.0				
3	K2	4	1 d	4	K1	2.0	b	1.0				
				5	K2	2.0	b	2.0				
	rig	ht		6	K4	2.0	b	3.0				
	k		rv	7	K1	3.0	С	1.0				
				- 8	K2	3.0	С	2.0				
	0	K1	1	. 9	K4	3.0	С	3.0				
	1	K2	2	10	K1	NaN	d	1.0				
	2	K4	3	11	K2	4.0	d	2.0				
				12	K4	4.0	d	3.0				

Another fairly common situation is to have two like-indexed (or similarly indexed) Series or DataFrame objects and wanting to "patch" values in one object from values for matching indices in the other. Here is an example:

For this, use the combine\_first method:

```
In [90]: result = df1.combine_first(df2)
```



Note that this method only takes values from the right DataFrame if they are missing in the left DataFrame. A related method, update, alters non-NA values inplace:

```
In [91]: df1.update(df2)
```

	df	f1			dt	f2		Result			
	0	1	2		0	1	2		0	1	2
0	NaN	3.0	5.0			_	-	0	NaN	3.0	5.0
1	-4.6	NaN	NaN	1	-42.6	NaN	-8.2	1	-42.6	NaN	-8.2
	4.0	TWOTE	TWOTE	2	-5.0	1.6	4.0		42.0	14014	-0.2
2	NaN	7.0	NaN					2	-5.0	1.6	4.0