

CHRIS ALBON

Learning machine learning? Try my [machine learning flashcards](#) or [Machine Learning with Python Cookbook](#).

Join And Merge Pandas Dataframe

20 Dec 2017

import modules

```
import pandas as pd
from IPython.display import display
from IPython.display import Image
```

Create a dataframe

```
raw_data = {
    'subject_id': ['1', '2', '3', '4', '5'],
    'first_name': ['Alex', 'Amy', 'Allen', 'Alice', 'Ayoung'],
    'last_name': ['Anderson', 'Ackerman', 'Ali', 'Aoni', 'Atiches']}
df_a = pd.DataFrame(raw_data, columns = ['subject_id', 'first_name',
    'last_name'])
df_a
```

	subject_id	first_name	last_name
0	1	Alex	Anderson
1	2	Amy	Ackerman
2	3	Allen	Ali
3	4	Alice	Aoni
4	5	Ayoung	Atiches

Create a second dataframe

CHRIS ALBON

```
    'first_name': ['Billy', 'Brian', 'Bran', 'Bryce', 'Betty'],  
    'last_name': ['Bonder', 'Black', 'Balwner', 'Brice', 'Btisan']}]  
df_b = pd.DataFrame(raw_data, columns = ['subject_id', 'first_name',  
    'last_name'])  
df_b
```

	subject_id	first_name	last_name
0	4	Billy	Bonder
1	5	Brian	Black
2	6	Bran	Balwner
3	7	Bryce	Brice
4	8	Betty	Btisan

Create a third dataframe

```
raw_data = {  
    'subject_id': ['1', '2', '3', '4', '5', '7', '8', '9', '10',  
    '11'],  
    'test_id': [51, 15, 15, 61, 16, 14, 15, 1, 61, 16]}  
df_n = pd.DataFrame(raw_data, columns = ['subject_id', 'test_id'])  
df_n
```

	subject_id	test_id
0	1	51
1	2	15
2	3	15
3	4	61
4	5	16
5	7	14
6	8	15
7	9	1
8	10	61
9	11	16

Join the two dataframes along rows

CHRIS ALBON

	subject_id	first_name	last_name
0	1	Alex	Anderson
1	2	Amy	Ackerman
2	3	Allen	Ali
3	4	Alice	Aoni
4	5	Ayoung	Atiches
0	4	Billy	Bonder
1	5	Brian	Black
2	6	Bran	Balwner
3	7	Bryce	Brice
4	8	Betty	Btisan

Join the two dataframes along columns

```
pd.concat([df_a, df_b], axis=1)
```

	subject_id	first_name	last_name	subject_id	first_name	last_name
0	1	Alex	Anderson	4	Billy	Bonder
1	2	Amy	Ackerman	5	Brian	Black
2	3	Allen	Ali	6	Bran	Balwner
3	4	Alice	Aoni	7	Bryce	Brice
4	5	Ayoung	Atiches	8	Betty	Btisan

Merge two dataframes along the subject_id value

```
pd.merge(df_new, df_n, on='subject_id')
```

	subject_id	first_name	last_name	test_id
0	1	Alex	Anderson	51
1	2	Amy	Ackerman	15
2	3	Allen	Ali	15

CHRIS ALBON

4	4	Billy	Bonder	61
5	5	Ayoung	Atiches	16
6	5	Brian	Black	16
7	7	Bryce	Brice	14
8	8	Betty	Btisan	15

Merge two dataframes with both the left and right dataframes using the `subject_id` key

```
pd.merge(df_new, df_n, left_on='subject_id', right_on='subject_id')
```

	subject_id	first_name	last_name	test_id
0	1	Alex	Anderson	51
1	2	Amy	Ackerman	15
2	3	Allen	Ali	15
3	4	Alice	Aoni	61
4	4	Billy	Bonder	61
5	5	Ayoung	Atiches	16
6	5	Brian	Black	16
7	7	Bryce	Brice	14
8	8	Betty	Btisan	15

Merge with outer join

“Full outer join produces the set of all records in Table A and Table B, with matching records from both sides where available. If there is no match, the missing side will contain null.” - [source](#)

```
pd.merge(df_a, df_b, on='subject_id', how='outer')
```

	subject_id	first_name_x	last_name_x	first_name_y	last_name_y
0	1	Alex	Anderson	NaN	NaN
1	2	Amy	Ackerman	NaN	NaN

CHRIS ALBON

3	4	Alice	Aoni	Billy	Bonder
4	5	Ayoung	Atiches	Brian	Black
5	6	NaN	NaN	Bran	Balwner
6	7	NaN	NaN	Bryce	Brice
7	8	NaN	NaN	Betty	Btisan

Merge with inner join

“Inner join produces only the set of records that match in both Table A and Table B.” - [source](#)

```
pd.merge(df_a, df_b, on='subject_id', how='inner')
```

	subject_id	first_name_x	last_name_x	first_name_y	last_name_y
0	4	Alice	Aoni	Billy	Bonder
1	5	Ayoung	Atiches	Brian	Black

Merge with right join

```
pd.merge(df_a, df_b, on='subject_id', how='right')
```

	subject_id	first_name_x	last_name_x	first_name_y	last_name_y
0	4	Alice	Aoni	Billy	Bonder
1	5	Ayoung	Atiches	Brian	Black
2	6	NaN	NaN	Bran	Balwner
3	7	NaN	NaN	Bryce	Brice
4	8	NaN	NaN	Betty	Btisan

Merge with left join

“Left outer join produces a complete set of records from Table A, with the matching records (where available) in Table B. If there is no match, the right side will contain null.” - [source](#)

```
pd.merge(df_a, df_b, on='subject_id', how='left')
```

CHRIS ALBON

1	2	Amy	Ackerman	NaN	NaN
2	3	Allen	Ali	NaN	NaN
3	4	Alice	Aoni	Billy	Bonder
4	5	Ayoung	Atiches	Brian	Black

Merge while adding a suffix to duplicate column names

```
pd.merge(df_a, df_b, on='subject_id', how='left', suffixes=('_left', '_right'))
```

	subject_id	first_name_left	last_name_left	first_name_right	last_name_right
0	1	Alex	Anderson	NaN	NaN
1	2	Amy	Ackerman	NaN	NaN
2	3	Allen	Ali	NaN	NaN
3	4	Alice	Aoni	Billy	Bonder
4	5	Ayoung	Atiches	Brian	Black

Merge based on indexes

```
pd.merge(df_a, df_b, right_index=True, left_index=True)
```

	subject_id_x	first_name_x	last_name_x	subject_id_y	first_name_y	last_name_y
0	1	Alex	Anderson	4	Billy	Bonder
1	2	Amy	Ackerman	5	Brian	Black
2	3	Allen	Ali	6	Bran	Balwner
3	4	Alice	Aoni	7	Bryce	Brice
4	5	Ayoung	Atiches	8	Betty	Btisan

Find an error or bug?

Everything on this site is available on GitHub. Head to [and submit a suggested change](#). You can also message me directly on [Twitter](#).

CHRIS ALBON