

Class_8_Capstone



default ▼

```
%pyspark
import pandas as pd
import numpy as np
```

FINISHED ▶ ↗↘ 📖 ⚙️

Took 45 sec. Last updated by anonymous at March 09 2017, 11:29:34 PM.

```
%pyspark
df = pd.DataFrame({'Key1':['a','a','b','b','a'],
                    'Key2':['one','two','one','two','one'],
                    'data1':np.random.randn(5),
                    'data2':np.random.randn(5)}
})
```

FINISHED ▶ ⌵ ⌶ ⚙

Took 28 sec. Last updated by anonymous at March 09 2017, 11:29:34 PM.

```
%pyspark
df
```

FINISHED ▶ ⌵ ⌶ ⚙

	Key1	Key2	data1	data2
0	a	one	-0.685101	-1.112390
1	a	two	0.951374	1.871486
2	b	one	1.337829	1.187213
3	b	two	0.005543	-0.820762
4	a	one	-0.331941	0.084290

Took 0 sec. Last updated by anonymous at March 09 2017, 11:29:39 PM.

```
%pyspark
grouped = df['data1'].groupby(df['Key1'])
```

FINISHED ▶ 🔍 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:29:43 PM.

```
%pyspark
grouped
```

FINISHED ▶ ↻ 📖 ⚙️

```
<pandas.core.groupby.SeriesGroupBy object at 0x00000000600F400>
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:29:45 PM.

```
%pyspark
grouped.mean()
```

FINISHED ▶ 🔍 📖 ⚙️

```
Key1
a    -0.021889
b     0.671686
Name: data1, dtype: float64
```

Took 1 sec. Last updated by anonymous at March 09 2017, 11:29:49 PM.

Class_8_Capstone

Untitled Untitled Untitled Untitled Untitled Untitled Untitled Untitled Untitled Untitled

Zeppelin

```
%pyspark
means = df['data1'].groupby([df['Key1'],df['Key2']]).mean()
```

Class_8_Capstone

FINISHED ▶ ⌵ 📖 ⚙️



default ▼

Took 0 sec. Last updated by anonymous at March 09 2017, 11:29:52 PM.

```
%pyspark
means
```

FINISHED ▶ ⌵ 📖 ⚙️

```
Key1  Key2
a      one   -0.508521
      two    0.951374
b      one    1.337829
      two    0.005543
Name: data1, dtype: float64
```

Took 1 sec. Last updated by anonymous at March 09 2017, 11:29:56 PM.

```
%pyspark
means.unstack()
```

FINISHED ▶ ⌵ 📖 ⚙️

```
Key2      one      two
Key1
a    -0.508521  0.951374
b     1.337829  0.005543
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:29:59 PM.

```
%pyspark
states = np.array(['ohio','california','california','ohio','ohio'])
```

FINISHED ▶ ⌵ 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:02 PM.

```
%pyspark
years = np.array([2005,2005,2006,2005,2006])
```

FINISHED ▶ ⌵ 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:07 PM.

```
%pyspark
df['data1'].groupby([states,years]).mean()
```

FINISHED ▶ ⌵ 📖 ⚙️

```
california 2005    0.951374
           2006    1.337829
ohio       2005   -0.339779
           2006   -0.331941
Name: data1, dtype: float64
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:09 PM.

```
%pyspark
df.groupby('Key1').mean()
```

FINISHED ▶ ⌵ 📖 ⚙️

Zeppelin

Key1

a -0.021889 0.281129

class 8 Can

Class_8_Capstone



default ▼

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:14 PM.

FINISHED ▶ 🔍 📖 ⚙️

```
%pyspark
df.groupby(['Key1', 'Key2']).mean()
```

		data1	data2
Key1	Key2		
a	one	-0.508521	-0.514050
	two	0.951374	1.871486
b	one	1.337829	1.187213
	two	0.005543	-0.820762

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:16 PM.

FINISHED ▶ ↗ ↘ 📖 ⚙️

```
%pyspark
df.groupby(['Key1', 'Key2']).size()
```

Key1	Key2	
a	one	2
	two	1
b	one	1
	two	1

```
dtype: int64
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:19 PM.

FINISHED ▶ ⌵ 📖 ⚙️

```
%pyspark
for name,group in df.groupby('Key1'):
    print name
    print group
```

```
a
  Key1 Key2    data1    data2
0    a  one -0.685101 -1.112390
1    a  two  0.951374  1.871486
4    a  one -0.331941  0.084290
b
  Key1 Key2    data1    data2
2    b  one  1.337829  1.187213
3    b  two  0.005543 -0.820762
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:23 PM.

FINISHED ▶ 🔍 📖 ⚙️

```
%pyspark
for (K1,K2),group in df.groupby(['Key1','Key2']):
    print K1,K2
    print group
```

0000 0 Constant

Key1

Key2

data1

data2

a two

1 a two 0.951374 1.871486

Key1

Key2

data1

data2

b one

2 b one 1.337829 1.187213

Key1

Key2

data1

data2

b two

3 b two 0.005543 -0.820762

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:27 PM.

▶

⌕

📖

✍

📄

⬇

📄

🗑

🕒

⌨

⚙

🔒

default ▾

%pyspark

pieces = dict(list(df.groupby('Key1')))

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:30 PM.

FINISHED ▶ ⌕ 📖 ⚙

%pyspark

pieces['b']

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:34 PM.

FINISHED ▶ ⌕ 📖 ⚙

Key1

Key2

data1

data2

2 b one 1.337829 1.187213

3 b two 0.005543 -0.820762

%pyspark

df.dtypes

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:37 PM.

FINISHED ▶ ⌕ 📖 ⚙

Key1

object

Key2

object

data1

float64

data2

float64

dtype:

object

%pyspark

grouped = df.groupby(df.dtypes,axis=1)

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:40 PM.

FINISHED ▶ ⌕ 📖 ⚙

%pyspark

dict(list(grouped))

Took 0 sec. Last updated by anonymous at March 09 2017, 11:30:40 PM.

FINISHED ▶ ⌕ 📖 ⚙

{dtype('O'):

Key1

Key2

0 a one

1 a two

2 b one

3 b two

4 a one, dtype('float64'):

data1

data2

Class_8_Capstone



FINISHED ▶ ⌵ ⌶ ⚙

FINISHED ▶ ⌵ ⌶ ⚙

FINISHED ▶ ⌵ ⌶ ⚙

FINISHED ▶ ⌵ ⌶ ⚙

FINISHED ▶ ⌵ ⌶ ⚙

FINISHED ▶ ⌵ ⌶ ⚙

Class 0 Constant

Class_8_Capstone

Took 0 sec. Last updated by anonymous at March 09 2017, 11:33:42 PM.

FINISHED ▶ ↗↘ 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:34:21 PM.

FINISHED ▶ ↻ 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:35:42 PM.

FINISHED ▶ ⌵ ⌶ ⚙

Took 0 sec. Last updated by anonymous at March 09 2017, 11:36:27 PM.

FINISHED ▶ ↻ 📖 ⚙️

a	red
b	red
c	blue

Class_8_Capstone

Class_o_Capstone Took 0 sec. Last updated by anonymous at March 09 2017, 11:37:04 PM

FINISHED ▶ ⌵ ⌶ ⚙

Took 0 sec. Last updated by anonymous at March 09 2017, 11:37:46 PM.

FINISHED ▶ ⌵ ⌶ ⚙

Took 0 sec. Last updated by anonymous at March 09 2017, 11:38:26 PM.

FINISHED ▶ 🔍 📖 ⚙️

Took 0 sec. Last updated by anonymous at March 09 2017, 11:39:20 PM.

FINISHED

Took 0 sec. Last updated by anonymous at March 09 2017, 11:39:55 PM.

FINISHED ▶ ↻ 📖 ⚙️

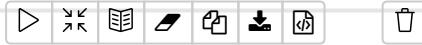
cty	US			JP		
tenor	1	3	5	1	3	
0	0.608612	-0.527968	1.540790	0.648186	-1.020264	
1	0.000860	-1.005398	-0.698190	-1.419141	0.540718	

2 0.063638 -3.539274 0.333320 -0.109230 0.823649
-0.042353 0.901137 -0.028447 -1.355646 0.506471

Zeppelin

Took 0 sec. Last updated by anonymous at March 09 2017, 11:40:37 PM.

Class_8_Capstone



```
%pyspark
hier_df.groupby(level='cty', axis=1).count()
```

FINISHED ▶ ⌵ 📖 ⚙

```
cty  JP  US
0     2   3
1     2   3
2     2   3
3     2   3
```

Took 0 sec. Last updated by anonymous at March 09 2017, 11:41:22 PM.

|

READY ▶ ⌵ 📖 ⚙