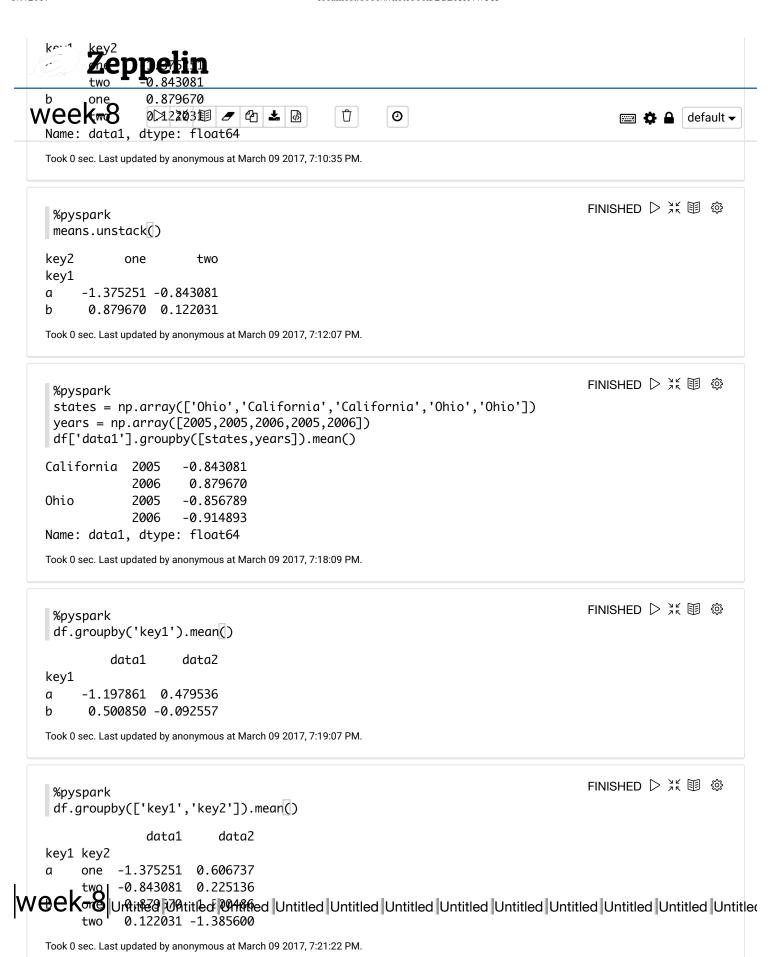


```
week-8
                  Û
                                                     ②
                                                                                                 default ▼
                                                                                  FINISHED ▷ 光 圓 ۞
   %pyspark
   import pandas as pd
   import numpy as np
   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)})
   df
        data1
                   data2 key1 key2
  0 -1.835609 0.600007
                            а
                                one
  1 -0.843081 0.225136
                                two
     0.879670 1.200486
                                one
  3 0.122031 -1.385600
                            b
                                two
  4 -0.914893 0.613466
                            а
                                one
  Took 0 sec. Last updated by anonymous at March 09 2017, 7:04:36 PM. (outdated)
                                                                                  FINISHED ▷ ※ ■ �
   %pyspark
   grouped = df['data1'].groupby(df['key1'])
   df
        data1
                   data2 key1 key2
  0 -1.835609 0.600007
                            а
                                one
  1 -0.843081 0.225136
                                two
  2 0.879670 1.200486
                                one
  3 0.122031 -1.385600
                            b
                                two
  4 -0.914893 0.613466
                            а
                                one
  Took 0 sec. Last updated by anonymous at March 09 2017, 7:06:20 PM.
                                                                                  FINISHED ▷ 💥 🗐 🕸
   %pyspark
  grouped.mean()
  key1
      -1.197861
       0.500850
  Name: data1, dtype: float64
  Took 0 sec. Last updated by anonymous at March 09 2017, 7:06:43 PM.
```

Untitled Unt



```
Zeppelin
                                                                               FINISHED ▷ 💥 🗉 🕸
 df.groupby(['key1','key2']).size()
               Ů
                                                  ②
                                                                                              default ▼
              2
      two
              1
b
      one
              1
      two
              1
dtype: int64
Took 0 sec. Last updated by anonymous at March 09 2017, 7:22:15 PM.
                                                                               FINISHED ▷ 端 圓 墩
 %pyspark
 for name, group in df.groupby('key1'):
      print name
      print group
а
      data1
                data2 key1 key2
0 -1.835609 0.600007
                             one
1 -0.843081 0.225136
                             two
4 -0.914893
             0.613466
                          а
                             one
b
      data1
                data2 key1 key2
2
  0.879670 1.200486
                          b
                             one
  0.122031 -1.385600
                          b
                            two
Took 0 sec. Last updated by anonymous at March 09 2017, 7:25:08 PM.
                                                                               FINISHED ▷ 牂 圓 ��
 %pyspark
 for (k1,k2), group in df.groupby(['key1','key2']):
    print k1, k2
    print group
a one
      data1
                data2 key1 key2
0 -1.835609 0.600007
                             one
4 -0.914893 0.613466
                             one
                          а
a two
      data1
                data2 key1 key2
1 -0.843081 0.225136
                          a two
b one
               data2 key1 key2
     data1
2 0.87967
           1.200486
                         b one
b two
      data1
              data2 key1 key2
3 0.122031 -1.3856
                        b two
Took 0 sec. Last updated by anonymous at March 09 2017, 7:27:50 PM.
```

