

Spark Python API函数学习 : pyspark API(1)

[《Spark Python API函数学习 : pyspark API\(1\)》](#)

[《Spark Python API函数学习 : pyspark API\(2\)》](#)

[《Spark Python API函数学习 : pyspark API\(3\)》](#)

[《Spark Python API函数学习 : pyspark API\(4\)》](#)

Spark支持Scala、Java以及Python语言，本文将通过图片和简单例子来学习pyspark API。



微信扫一扫，加关注
即可及时了解Spark、Hadoop或者Hbase
等相关的文章
欢迎关注微信公共帐号: iteblog_hadoop

过往记忆博客 (<http://www.iteblog.com>)
专注于Hadoop、Spark、Flume、Hbase等
技术的博客，欢迎关注。

Hadoop、Hive、Hbase、Flume等交流群：138615359和149892483

如果想及时了
解Spark、Hadoop或者Hbase相关的文章，欢迎关注微信公共帐号 : iteblog_hadoop

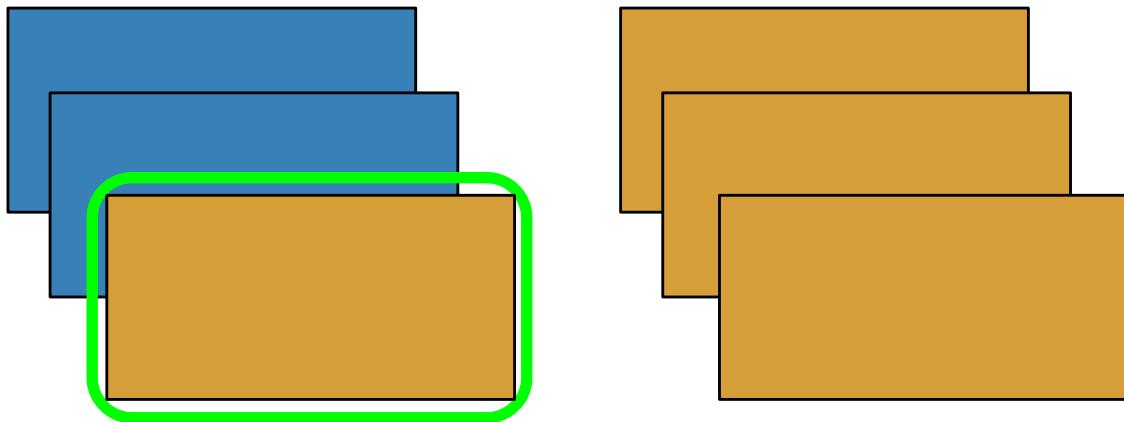
pyspark version

```
# print Spark version
print("pyspark version:" + str(sc.version))
```

pyspark version:1.2.2

map

—

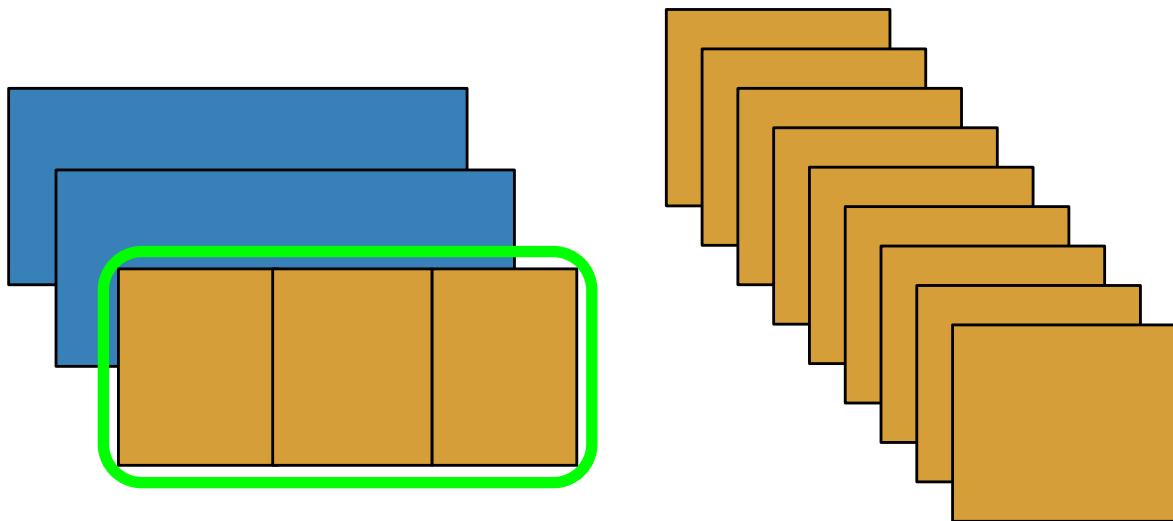


```
# map
# sc = spark context, parallelize creates an RDD from the passed object
x = sc.parallelize([1,2,3])
y = x.map(lambda x: (x,x**2))

# collect copies RDD elements to a list on the driver
print(x.collect())
print(y.collect())

[1, 2, 3]
[(1, 1), (2, 4), (3, 9)]
```

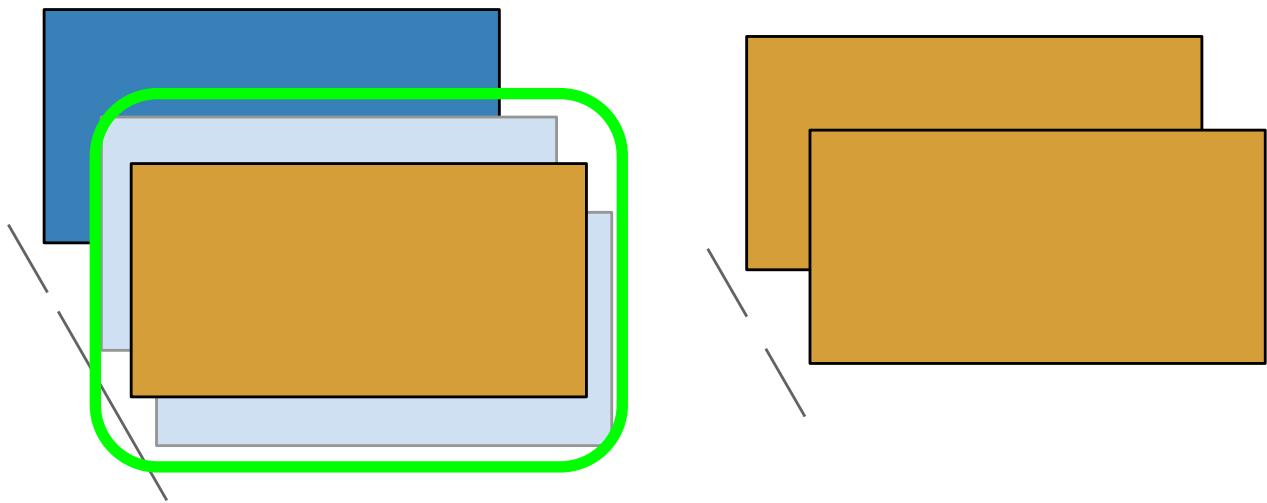
flatMap



```
# flatMap
x = sc.parallelize([1,2,3])
y = x.flatMap(lambda x: (x, 100*x, x**2))
print(x.collect())
print(y.collect())
```

```
[1, 2, 3]
[1, 100, 1, 2, 200, 4, 3, 300, 9]
```

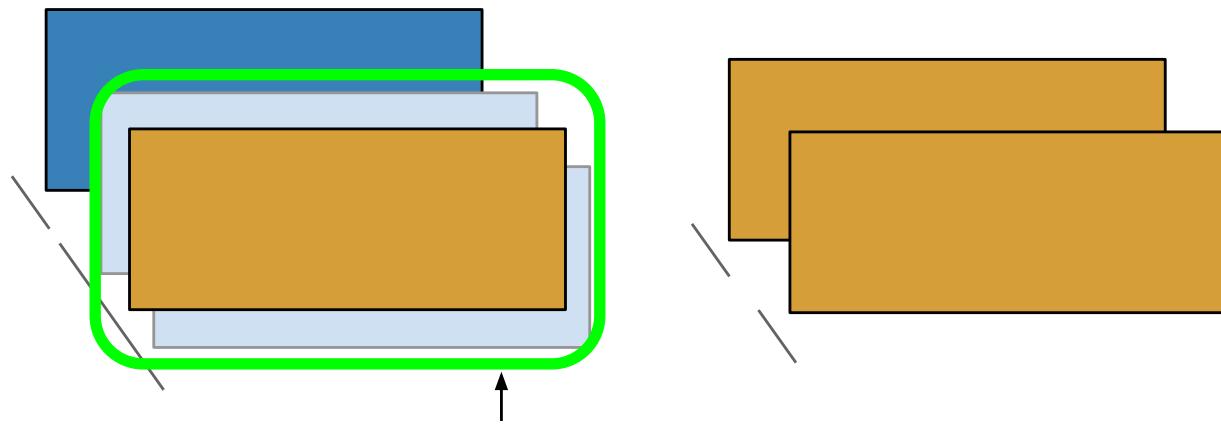
mapPartitions



```
# mapPartitions
x = sc.parallelize([1,2,3], 2)
def f(iterator): yield sum(iterator)
y = x.mapPartitions(f)
# glom() flattens elements on the same partition
print(x.glom().collect())
print(y.glom().collect())
```

```
[[1], [2, 3]]
[[1], [5]]
```

mapPartitionsWithIndex



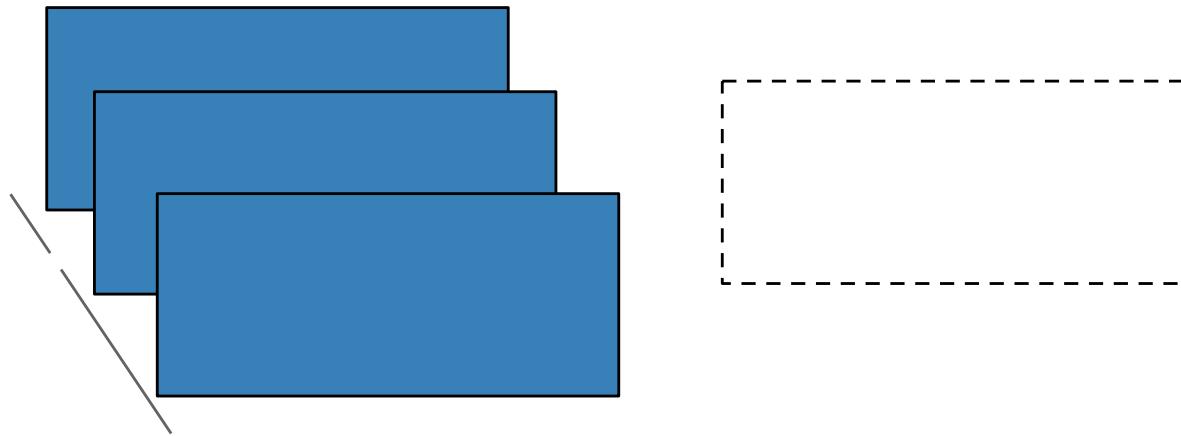
```
# mapPartitionsWithIndex
x = sc.parallelize([1,2,3], 2)
def f(partitionIndex, iterator): yield (partitionIndex,sum(iterator))
y = x.mapPartitionsWithIndex(f)
```

```
# glom() flattens elements on the same partition
print(x.glom().collect())
print(y.glom().collect())
```

```
[[1], [2, 3]]
[(0, 1), (1, 5)]
```

getNumPartitions

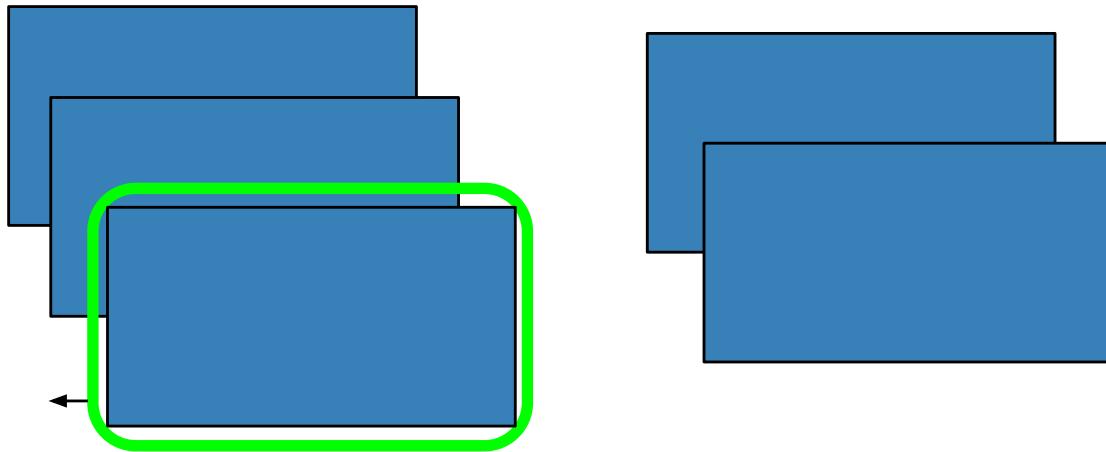
2



```
#getNumPartitions
x = sc.parallelize([1,2,3], 2)
y = x.getNumPartitions()
print(x.glom().collect())
print(y)
```

```
[[1], [2, 3]]
2
```

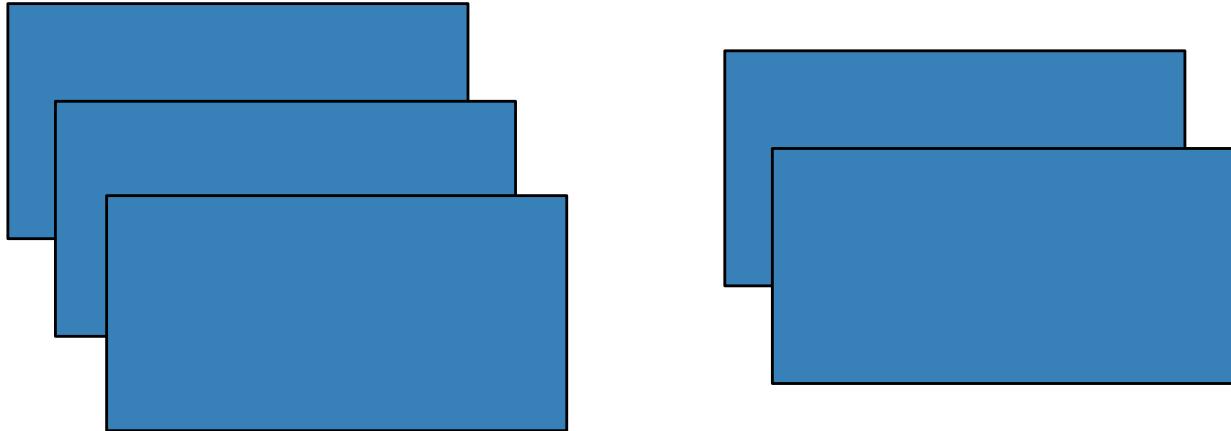
filter



```
# filter
x = sc.parallelize([1,2,3])
y = x.filter(lambda x: x%2 == 1) # filters out even elements
print(x.collect())
print(y.collect())
```

```
[1, 2, 3]
[1, 3]
```

distinct

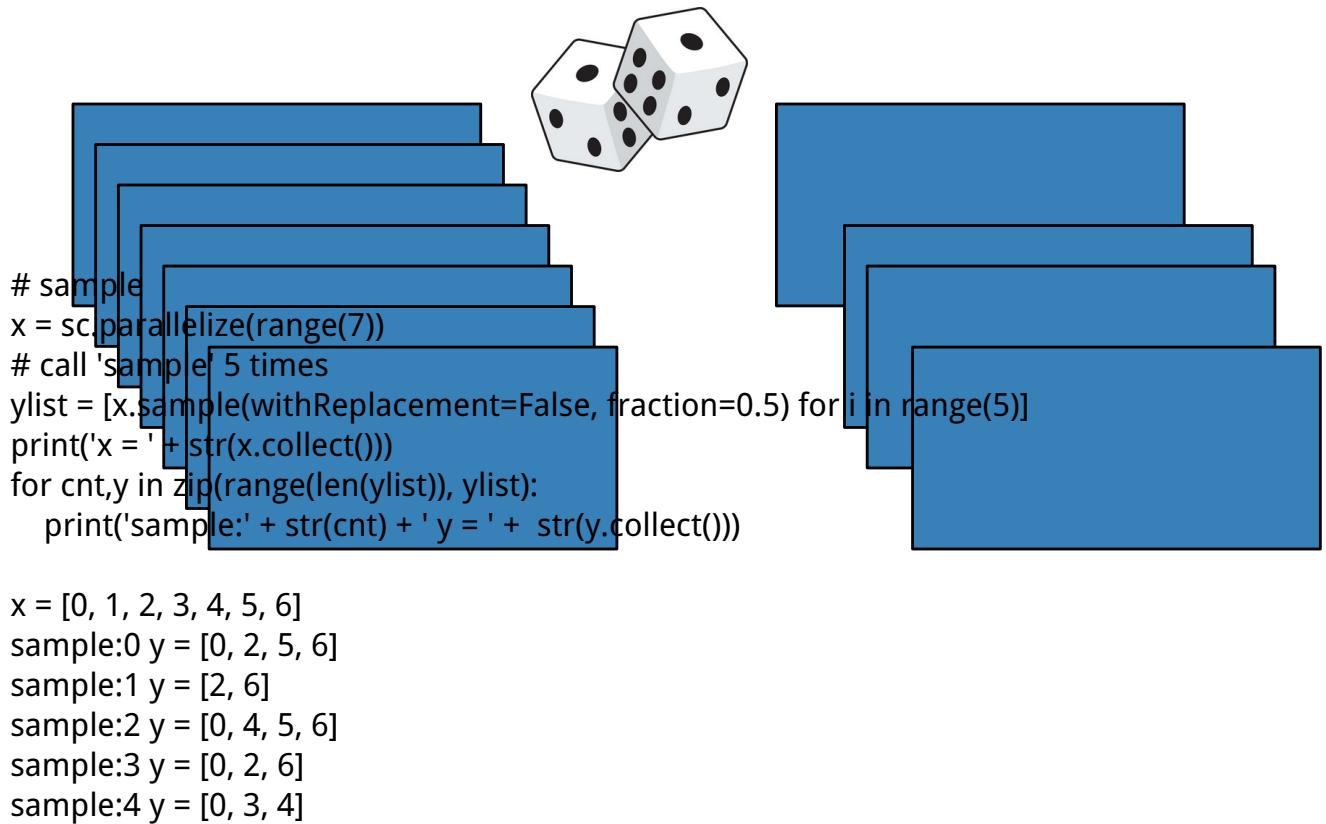


```
# distinct
x = sc.parallelize(['A','A','B'])
y = x.distinct()
print(x.collect())
print(y.collect())
```

```
['A', 'A', 'B']
['A', 'B']
```

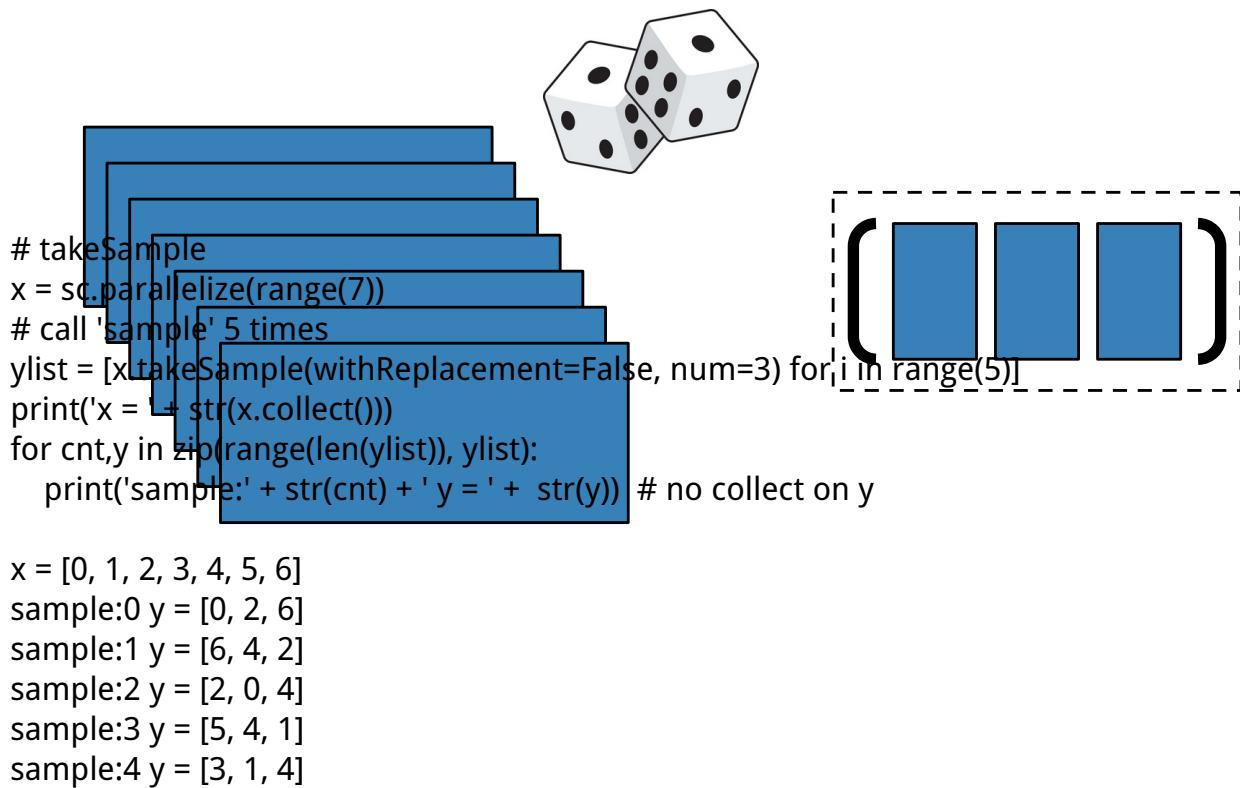
sample

■

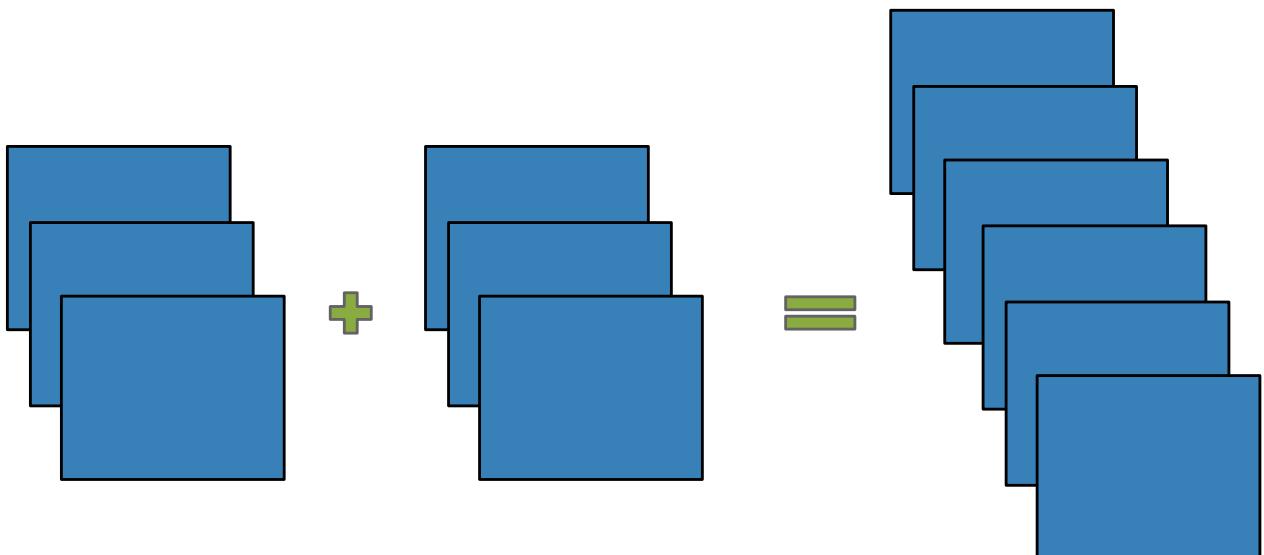


takeSample

—



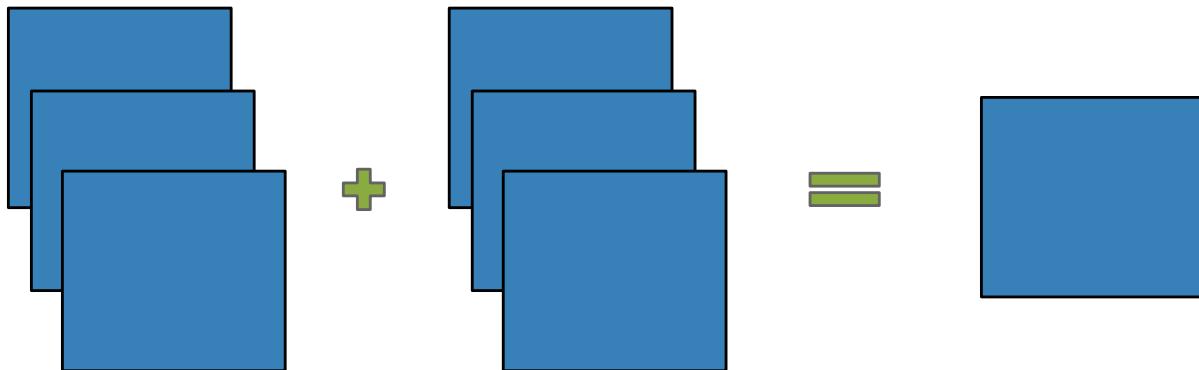
union



```
# union
x = sc.parallelize(['A','A','B'])
y = sc.parallelize(['D','C','A'])
z = x.union(y)
print(x.collect())
print(y.collect())
print(z.collect())

['A', 'A', 'B']
['D', 'C', 'A']
['A', 'A', 'B', 'D', 'C', 'A']
```

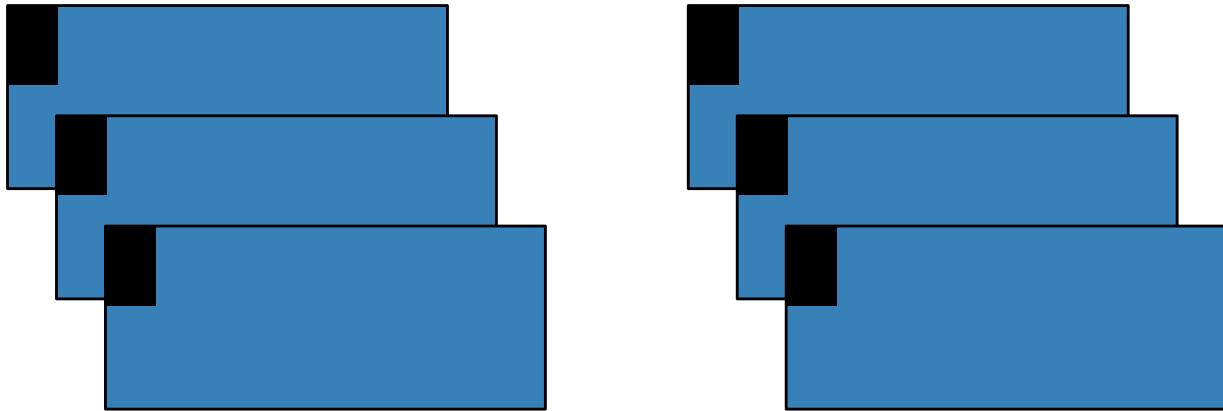
intersection



```
# intersection
x = sc.parallelize(['A','A','B'])
y = sc.parallelize(['A','C','D'])
z = x.intersection(y)
print(x.collect())
print(y.collect())
print(z.collect())
```

```
['A', 'A', 'B']
['A', 'C', 'D']
['A']
```

sortByKey



```
# sortByKey
x = sc.parallelize([('B',1),('A',2),('C',3)])
y = x.sortByKey()
print(x.collect())
print(y.collect())
```

```
[('B', 1), ('A', 2), ('C', 3)]
[('A', 2), ('B', 1), ('C', 3)]
```

本博客文章除特别声明，全部都是原创！

原创文章版权归过往记忆大数据（过往记忆）所有，未经许可不得转载。

本文链接: 【】()