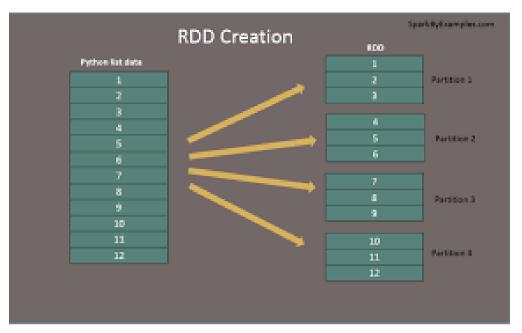
```
1 df 4.show()
In [139]:
                 Name | Departments | salary |
                Krish|Data Science| 10000|
                Krishl
                                IOT|
                                      50001
               Maheshl
                          Big Datal
                                      4000
                          Big Data
               Krishl
                                      4000
               Mahesh|Data Science|
                                      3000
           |Sudhanshu|Data Science|
                                     20000
           |Sudhanshu|
                                IOT| 10000|
           | Sudhanshu|
                          Big Data|
                                      50001
                Sunny Data Science | 10000 |
                Sunny|
                          Big Data | 2000 |
```

## RDD

In [ ]:

RDD (Resilient Distributed Dataset) is a fundamental building block of Pyspark which is fault tolerant, immutable distributed collection of objects.

Once you create a record in an RDD you cannot change it. Each record in RDD is divided into logical partition which can be computed on different node of cluster.



## creatting rdd

ParallelCollectionRDD[28] at readRDDFromFile at PythonRDD.scala:274

```
In [13]: 1 type(rdd1)
```

Out[13]: pyspark.rdd.RDD

```
1 rdd1.collect()
In [38]:
Out[38]: [23, 45, 67, 86, 78, 3, 4, 5, 6, 10, 11, 12, 23, 45, 67, 10]
          1 rdd1.count()
In [39]:
Out[39]: 16
          1 type(result)
In [18]:
Out[18]: list
In [21]:
             for val in rdd1.collect():
          2
                 print(val)
         23
         45
         67
         86
         78
         10
         11
         12
In [23]:
             rdd1.count()
Out[23]: 12
In [35]:
          1 rdd1.distinct().collect()
Out[35]: [10, 67, 3, 11, 4, 12, 45, 5, 86, 78, 6, 23]
          1 rdd1.distinct().count()
In [40]:
Out[40]: 12
```

```
In [ ]:
          1
          1 rdd1.collect()
In [41]:
Out[41]: [23, 45, 67, 86, 78, 3, 4, 5, 6, 10, 11, 12, 23, 45, 67, 10]
          1 rdd1.filter(lambda x : x<=20).collect()</pre>
In [42]:
Out[42]: [3, 4, 5, 6, 10, 11, 12, 10]
In [43]:
          1 rdd1.collect()
Out[43]: [23, 45, 67, 86, 78, 3, 4, 5, 6, 10, 11, 12, 23, 45, 67, 10]
In [ ]:
In [44]:
          1 rdd1.first()
Out[44]: 23
In [45]:
          1 rdd1.take(5)
Out[45]: [23, 45, 67, 86, 78]
In [ ]:
In [46]:
          1 rdd1.collect()
Out[46]: [23, 45, 67, 86, 78, 3, 4, 5, 6, 10, 11, 12, 23, 45, 67, 10]
In [47]:
          1 rdd1.reduce(lambda x,y : x+y)
Out[47]: 495
          1 rdd1.saveAsTextFile('file.txt')
In [48]:
```

```
In [ ]:
          1
          1 rdd2 = spark.sparkContext.parallelize([1,2,3,4,5])
In [49]:
In [50]:
          1 rdd2
Out[50]: ParallelCollectionRDD[43] at readRDDFromFile at PythonRDD.scala:274
In [51]:
          1 rdd2.map(lambda x: x**3).collect()
Out[51]: [1, 8, 27, 64, 125]
In [ ]:
In [ ]:
          1 rdd3 = spark.sparkContext.parallelize([2,4,5,6,7,8,9])
In [52]:
          1 uni1 = rdd3.filter(lambda x : x%2 == 0)
In [54]:
          2 uni2 = rdd3.filter(lambda x : x%3 == 0)
In [55]:
          1 uni1.collect()
Out[55]: [2, 4, 6, 8]
          1 uni2.collect()
In [56]:
Out[56]: [6, 9]
          1 final = unil.union(uni2)
In [58]:
          1 final.collect()
In [59]:
Out[59]: [2, 4, 6, 8, 6, 9]
```

```
1 from pyspark.mllib.linalg import Matrix, Matrices
In [60]:
In [62]:
             data = [10,20,30,40,50,60]
           3 res = Matrices.dense(3,2,data)
In [63]:
           1 res
Out[63]: DenseMatrix(3, 2, [10.0, 20.0, 30.0, 40.0, 50.0, 60.0], False)
In [64]:
          1 print(res)
         DenseMatrix([[10., 40.],
                       [20., 50.],
                       [30., 60.]])
In [65]:
          1 type(res)
Out[65]: pyspark.mllib.linalg.DenseMatrix
         it returns seperate value for each element in RDD -- flatmap
           1 data = spark.sparkContext.parallelize(["Hey There", "This is RDD Session in Pyspark"])
In [66]:
In [67]:
          1 data.collect()
Out[67]: ['Hey There', 'This is RDD Session in Pyspark']
          1 data.flatMap(lambda x:x.split(" ")).collect()
In [68]:
Out[68]: ['Hey', 'There', 'This', 'is', 'RDD', 'Session', 'in', 'Pyspark']
In [ ]:
           1
           1 marks = [('Punit',55),('Salam',70),('Dharmesh',80),('Rohan',80),('Amit',55),('Sumit',90),('Punit',60),('
In [73]:
```

```
In [70]:
           1 marks
Out[70]: [('Punit', 55),
          ('Salam', 70),
           ('Dharmesh', 80),
           ('Rohan', 80),
           ('Amit', 55),
           ('Sumit', 90)]
In [74]:
           1 rdd5 = spark.sparkContext.parallelize(marks)
           1 rdd5.collect()
In [75]:
Out[75]: [('Punit', 55),
          ('Salam', 70),
           ('Dharmesh', 80),
           ('Rohan', 80),
           ('Amit', 55),
           ('Sumit', 90),
           ('Punit', 60),
           ('Dharmesh', 70)1
```

## reduce by key:

It performs multiple parallel process for each key in the data and combines the value for same key

It uses lambda to perfrom task

```
In [ ]:
           1 rdd5.sortByKey('ascending').collect()
In [77]:
Out[77]: [('Amit', 55),
          ('Dharmesh', 80),
           ('Dharmesh', 70),
           ('Punit', 55),
           ('Punit', 60),
           ('Rohan', 80),
           ('Salam', 70),
           ('Sumit', 90)]
           1 rdd5.sortByKey(ascending=False).collect()
In [79]:
Out[79]: [('Sumit', 90),
           ('Salam', 70),
           ('Rohan', 80),
           ('Punit', 55),
           ('Punit', 60),
           ('Dharmesh', 80),
           ('Dharmesh', 70),
           ('Amit', 55)]
           1 rdd5.collect()
In [81]:
Out[81]: [('Punit', 55),
           ('Salam', 70),
           ('Dharmesh', 80),
           ('Rohan', 80),
           ('Amit', 55),
           ('Sumit', 90),
           ('Punit', 60),
           ('Dharmesh', 70)]
           1 result = rdd5.groupByKey().collect()
In [83]:
```

```
In [851:
           1 for key, val in result:
                  print(key,list(val))
           2
          Punit [55, 60]
          Dharmesh [80, 70]
          Sumit [90]
          Rohan [80]
          Amit [55]
          Salam [70]
In [86]:
           1 count = rdd5.countByKey().items()
In [87]:
           1 for k,v in count:
           2
                  print(k,v)
          Punit 2
          Salam 1
          Dharmesh 2
          Rohan 1
          Amit 1
          Sumit 1
 In [ ]:
           1
              empty = spark.sparkContext.emptyRDD
In [88]:
In [91]:
           1 empty
Out[91]: <bound method SparkContext.emptyRDD of <SparkContext master=local[*] appName=session1>>
           1 lst = [6,7,8,9,98,87,66,54,33,54]
In [98]:
           1 rdd6 = spark.sparkContext.parallelize(lst,5)
In [99]:
In [101]:
           1 rdd6.saveAsTextFile('demo.txt')
```

```
1 rdd6.getNumPartitions()
In [102]:
Out[102]: 5
           1 rdd6.max()
In [103]:
Out[103]: 98
           1 rdd6.min()
In [104]:
Out[104]: 6
In [106]:
           1 rdd6.collect()
Out[106]: [6, 7, 8, 9, 98, 87, 66, 54, 33, 54]
 In [ ]:
 In [ ]:
In [108]:
           1 rdd5.collect()
Out[108]: [('Punit', 55),
           ('Salam', 70),
           ('Dharmesh', 80),
           ('Rohan', 80),
           ('Amit', 55),
           ('Sumit', 90),
           ('Punit', 60),
           ('Dharmesh', 70)]
In [109]:
           1 cols = ["Name", "Marks"]
In [110]:
           1 marks df = rdd5.toDF(cols)
```

```
In [111]:
           1 marks_df.show()
               Name|Marks|
              Punit|
                       55|
              Salam
                       70 j
          | Dharmesh |
                       80|
              Rohan
                       80|
               Amit
                       55
              Sumit
                       90|
              Punit
                       60|
          Dharmesh
                       70|
          +----+
In [112]:
           1 col = ["Data"]
In [114]:
           1 from pyspark.sql import Row
           1 row = Row("Data")
In [115]:
In [117]:
           1 res = rdd6.map(row).toDF()
```

```
In [118]: 1 res.show()

----+
| Data|
----+
| 6|
| 7|
| 8|
| 9|
| 98|
| 87|
| 66|
| 54|
| 33|
| 54|
+---+
```