

```

1 package com.test.spark3
2
3 import org.apache.spark.storage.StorageLevel
4 import org.apache.spark.{HashPartitioner, SparkConf, SparkContext}
5
6 object visitCountProv {
7   def main(args: Array[String]): Unit = {
8     val conf = new SparkConf().setAppName("test").setMaster("local")
9     conf.set("spark.testing.memory", "471859200")
10    val sc = new SparkContext(conf)    //Driver类
11
12    val rdd = sc.textFile("D:\\doc\\study\\51CTO\\大数据中级\\day6-spark3\\
    电商流量数据文件\\data\\2015082818")
13    .cache()          //缓存在内存和HDFS的/tmp/Spark* 路径下
14    .filter(line => line.length>0)
15    .map { line=>
16      val arr = line.split("\\t")
17      val url = arr(1)
18      val guid = arr(5)
19      val provinceId = arr(23)
20      val date = arr(17).substring(0,10)    //2015-08-28 18:10:00
21      ((date,provinceId, guid),url)
22    }.filter( line=> line._2.length>5)    // length(url)>5
23    .partitionBy(new HashPartitioner(10))
24    .persist(StorageLevel.DISK_ONLY)
25
26    /*
27      select date, provinceId,count(url) pv ,count(distinct guid) uv
28      from track_log
29      group by date,provinceId
30    */
31
32    rdd.mapValues(url=>1)    // ( (date,provinceId, guid), 1)
33    .reduceByKey(_+_ )      // ((date,provinceId, guid), pv) pv=count(url)
34    .map{ case ((date, provId, guid),pv) => ((date,provId), (pv,1)) }    // ((date
,provinceId), (pv, uv))
35    .reduceByKey((x,y)=>(x._1+y._1, x._2+y._2))
36    .sortByKey(true)
37    .foreach(println)
38
39    rdd.unpersist()
40    sc.stop()
41  }
42 }
43

```