

深入理解Spark 2.1 Core（十）：Shuffle Map 端的原理与源码分析

■ 版权声明：本文为博主原创文章，转载请附上原文地址。 <http://blog.csdn.net/u011239443/article/details/55044862>

<http://blog.csdn.net/u011239443/article/details/55044862>

在上一篇《深入理解Spark 2.1 Core（九）：迭代计算和Shuffle的原理与源码分析》提到经过迭代计算后，`SortShuffleWriter.write`中：

```
1 // 根据排序方式，对数据进行排序并写入内存缓冲区。
2 // 若排序中计算结果超出的阈值，
3 // 则将其溢写到磁盘数据文件
4 sorter.insertAll(records)
```

我们先来宏观的了解下Map端，我们会根据 `aggregator.isDefined` 是否定义了聚合函数和 `ordering.isDefined` 是否定义了排序函数分为三种：

- 没有聚合和排序，数据先按照partition写入不同的文件中，最后按partition顺序合并写入同一文件。适合partition数量较少时。将多个bucket合并到同一文件，减少map输出文件数，节省磁盘I/O，提高性能。
- 没有聚合但有排序，在缓存对数据先根据分区（或者还有key）进行排序，最后按partition顺序合并写入同一文件。适合当partition数量较多时。将多个bucket合并到同一文件，减少map输出文件数，节省磁盘I/O，提高性能。缓存使用超过阈值，将数据写入磁盘。
- 有聚合有排序，现在缓存中根据key值聚合，再在缓存对数据先根据分区（或者还有key）进行排序，最后按partition顺序合并写入同一文件。将多个bucket合并到同一文件，减少map输出文件数，节省磁盘I/O，提高性能。缓存使用超过阈值，将数据写入磁盘。逐条的读取数据，并进行聚合，减少了内存的占用。

我们先来深入看下 `insertAll`：

```
1 def insertAll(records: Iterator[Product2[K, V]]): Unit = {
2   // 若定义了聚合函数，则shouldCombine为true
3   val shouldCombine = aggregator.isDefined
4
5   // 外部排序是否需要聚合
6   if (shouldCombine) {
7     // mergeValue 是对 Value 进行 merge 的函数
8     val mergeValue = aggregator.get.mergeValue
9     // createCombiner 为生成 Combiner 的函数
10    val createCombiner = aggregator.get.createCombiner
11    var kv: Product2[K, V] = null
12    // update 为偏函数
13    val update = (hadValue: Boolean, oldValue: C) => {
14      // 当有Value时，将oldValue与新的Value kv._2 进行merge
15      // 若没有Value，传入kv._2，生成Value
16      if (hadValue) mergeValue(oldValue, kv._2) else createCombiner(kv._2)
17    }
18    while (records.hasNext) {
19      addElementsRead()
20      kv = records.next()
21      // 首先使用我们的AppendOnlyMap
22      // 在内存中对value进行聚合
23      map.changeValue((getPartition(kv._1), kv._1), update)
24      // 超过阈值时写入磁盘
25      maybeSpillCollection(usingMap = true)
26    }
27  } else {
```

```

28      // 直接把Value插入缓冲区
29      while (records.hasNext) {
30          addElementsRead()
31          val kv = records.next()
32          buffer.insert(getPartition(kv._1), kv._1, kv._2.asInstanceOf[C])
33          maybeSpillCollection(usingMap = false)
34      }
35  }
36  }

```

这里的 `createCombiner` 我们可以看做用 `kv._2` 生成一个 `Value`。而 `mergeValue` 我们可以理解成为 `MapReduce` 中的 `combiner`，即可以理解成 Map 端的 Reduce 操作，先对相同的 `key` 的 `Value` 进行聚合。

聚合算法

下面我们来深入看看聚合操作部分：

调用栈：

- `util.collection.SizeTrackingAppendOnlyMap.changeValue`
 - `util.collection.AppendOnlyMap.changeValue`
 - `util.collection.AppendOnlyMap.incrementSize`
 - `util.collection.AppendOnlyMap.growTable`
 - `util.collection.SizeTracker.afterUpdate`
 - `util.collection.SizeTracker.takeSample`

首先是 `AppendOnlyMap` 的 `changeValue` 函数：

util.collection.SizeTrackingAppendOnlyMap.changeValue

```

1  override def changeValue(key: K, updateFunc: (Boolean, V) => V): V = {
2      // 应用聚合算法得到newValue
3      val newValue = super.changeValue(key, updateFunc)
4      // 更新对 AppendOnlyMap 大小的采样
5      super.afterUpdate()
6      // 返回结果
7      newValue
8  }

```

util.collection.AppendOnlyMap.changeValue

聚合算法：

```

1  def changeValue(key: K, updateFunc: (Boolean, V) => V): V = {
2      assert(!destroyed, destructionMessage)
3      val k = key.asInstanceOf[AnyRef]
4      if (k.eq(null)) {
5          if (!haveNullValue) {
6              incrementSize()
7          }
8          nullValue = updateFunc(haveNullValue, nullValue)
9          haveNullValue = true
10         return nullValue
11     }
12     // 根据k的hashCode在哈希 与 上 掩码 得到 pos
13     // 2*pos 为 k 应该所在的位置

```

```

14 // 2*pos + 1 为 k 对应的 v 所在的位置
15 var pos = rehash(k.hashCode) & mask
16 var i = 1
17 while (true) {
18 // 得到data中k所在的位置上的值curKey
19 val curKey = data(2 * pos)
20 if (curKey.eq(null)) {
21 // 若curKey为空
22 // 得到根据 kv._2, 即单个新值 生成的 newValue
23 val newValue = updateFunc(false, null.asInstanceOf[V])
24 data(2 * pos) = k
25 data(2 * pos + 1) = newValue.asInstanceOf[AnyRef]
26 // 扩充容量
27 incrementSize()
28 return newValue
29 } else if (k.eq(curKey) || k.equals(curKey)) {
30 // 若k 与 curKey 相等
31 // 将oldValue (data(2 * pos + 1)) 和 新的Value (kv._2) 进行聚合
32 val newValue = updateFunc(true, data(2 * pos + 1).asInstanceOf[V])
33 data(2 * pos + 1) = newValue.asInstanceOf[AnyRef]
34 return newValue
35 } else {
36 // 若curKey 不为null, 也和k不想等,
37 // 即 hash 冲突
38 // 则 不断的向后遍历 直到出现前两种情况
39 val delta = i
40 pos = (pos + delta) & mask
41 i += 1
42 }
43 }
44 null.asInstanceOf[V]
45 }

```

util.collection.AppendOnlyMap.incrementSize

我们再来看一下扩充容量的实现：

```

1 private def incrementSize() {
2   curSize += 1
3   // 当curSize大于阈值growThreshold时,
4   // 调用growTable()
5   if (curSize > growThreshold) {
6     growTable()
7   }
8 }

```

util.collection.AppendOnlyMap.growTable

```

1 protected def growTable() {
2   生成容量翻倍的新Data
3   val newCapacity = capacity * 2
4   require(newCapacity <= MAXIMUM_CAPACITY, s"Can't contain more than ${growThreshold} elements")
5   val newData = new Array[AnyRef](2 * newCapacity)
6   // 生成newMask
7   val newMask = newCapacity - 1
8   var oldPos = 0
9   while (oldPos < capacity) {
10    // 将旧的数据中的数据用newMask重新计算位置,
11    // 复制到新的Data 中
12    if (!data(2 * oldPos).eq(null)) {
13      val key = data(2 * oldPos)
14      val value = data(2 * oldPos + 1)
15      var newPos = rehash(key.hashCode) & newMask

```

```

16     var i = 1
17     var keepGoing = true
18     while (keepGoing) {
19         val curKey = newData(2 * newPos)
20         if (curKey.eq(null)) {
21             newData(2 * newPos) = key
22             newData(2 * newPos + 1) = value
23             keepGoing = false
24         } else {
25             val delta = i
26             newPos = (newPos + delta) & newMask
27             i += 1
28         }
29     }
30 }
31 oldPos += 1
32 }
33 // 更新
34 data = newData
35 capacity = newCapacity
36 mask = newMask
37 growThreshold = (LOAD_FACTOR * newCapacity).toInt
38 }

```

util.collection.SizeTracker.afterUpdate

我们回过头来看 `SizeTrackingAppendOnlyMap.changeValue` 中的更新对 `AppendOnlyMap` 大小的采样 `super.afterUpdate()`。所谓大小的采样，是只一次 `Update` 后 `AppendOnlyMap` 大小的变化量。但是如果在每次如 `insert`update` 等操作后就进行计算一次 `AppendOnlyMap` 会大大降低性能。所以，这里采用了采样估计的方法：

```

1     protected def afterUpdate(): Unit = {
2         numUpdates += 1
3         // 若numUpdates到达阈值，
4         // 则进行采样
5         if (nextSampleNum == numUpdates) {
6             takeSample()
7         }
8     }

```

util.collection.SizeTracker.takeSample

```

1     private def takeSample(): Unit = {
2         samples.enqueue(Sample(SizeEstimator.estimate(this), numUpdates))
3         // 只用两个采样
4         if (samples.size > 2) {
5             samples.dequeue()
6         }
7         val bytesDelta = samples.toList.reverse match {
8             // 估计出每次更新的变化量
9             case latest :: previous :: tail =>
10                 (latest.size - previous.size).toDouble / (latest.numUpdates - previous.numUpdates)
11             // 若小于 2个 样本，假设没产生变化
12             case _ => 0
13         }
14         // 更新
15         bytesPerUpdate = math.max(0, bytesDelta)
16         // 增大阈值
17         nextSampleNum = math.ceil(numUpdates * SAMPLE_GROWTH_RATE).toLong
18     }

```

我们再看来下估计 `AppendOnlyMap` 大小的函数：

```

1  def estimateSize(): Long = {
2      assert(samples.nonEmpty)
3      // 计算估计的总变化量
4      val extrapolatedDelta = bytesPerUpdate * (numUpdates - samples.last.numUpdates)
5      // 之前的大小 加上 估计的总变化量
6      (samples.last.size + extrapolatedDelta).toLong
7  }

```

写缓冲区

现在我们回到 `insertAll`，深入看看如何直接把 `Value` 插入缓冲区。

调用栈：

- `util.collection.PartitionedPairBuffer.insert`
 - `util.collection.PartitionedPairBuffer.growArray`

`util.collection.PartitionedPairBuffer.insert`

```

1  def insert(partition: Int, key: K, value: V): Unit = {
2      // 到了容量大小, 调用growArray()
3      if (curSize == capacity) {
4          growArray()
5      }
6      data(2 * curSize) = (partition, key.asInstanceOf[AnyRef])
7      data(2 * curSize + 1) = value.asInstanceOf[AnyRef]
8      curSize += 1
9      afterUpdate()
10 }

```

`util.collection.PartitionedPairBuffer.growArray`

```

1  private def growArray(): Unit = {
2      if (capacity >= MAXIMUM_CAPACITY) {
3          throw new IllegalStateException(s"Can't insert more than ${MAXIMUM_CAPACITY} elements")
4      }
5      val newCapacity =
6          if (capacity * 2 < 0 || capacity * 2 > MAXIMUM_CAPACITY) { // Overflow
7              MAXIMUM_CAPACITY
8          } else {
9              capacity * 2
10         }
11         // 生成翻倍容量的newArray
12         val newArray = new Array[AnyRef](2 * newCapacity)
13         // 复制
14         System.arraycopy(data, 0, newArray, 0, 2 * capacity)
15         data = newArray
16         capacity = newCapacity
17         resetSamples()
18     }

```

溢出

现在我们回到 `insertAll`，深入看看如何将超过阈值时写入磁盘：

调用栈：

- `util.collection.ExternalSorter.maybeSpillCollection`
 - `util.collection.Spillable.maybeSpill`

- util.collection.Spillable.spill
 - util.collection.ExternalSorter.spillMemoryIteratorToDisk

util.collection.ExternalSorter.maybeSpillCollection

```

1  private def maybeSpillCollection(usingMap: Boolean): Unit = {
2      var estimatedSize = 0L
3      if (usingMap) {
4          estimatedSize = map.estimateSize()
5          if (maybeSpill(map, estimatedSize)) {
6              map = new PartitionedAppendOnlyMap[K, C]
7          }
8      } else {
9          estimatedSize = buffer.estimateSize()
10         if (maybeSpill(buffer, estimatedSize)) {
11             buffer = new PartitionedPairBuffer[K, C]
12         }
13     }
14
15     if (estimatedSize > _peakMemoryUsedBytes) {
16         _peakMemoryUsedBytes = estimatedSize
17     }
18 }

```

util.collection.Spillable.maybeSpill

```

1  protected def maybeSpill(collection: C, currentMemory: Long): Boolean = {
2      var shouldSpill = false
3      if (elementsRead % 32 == 0 && currentMemory >= myMemoryThreshold) {
4          // 若大于阈值
5          // amountToRequest 为要申请的内存空间
6          val amountToRequest = 2 * currentMemory - myMemoryThreshold
7          val granted = acquireMemory(amountToRequest)
8          myMemoryThreshold += granted
9          // 如果我们分配了太小的内存,
10         // 由于 tryToAcquire 返回0
11         // 或者 内存申请大小超过了myMemoryThreshold
12         // 导致 依然 currentMemory >= myMemoryThreshold
13         // 则 shouldSpill
14         shouldSpill = currentMemory >= myMemoryThreshold
15     }
16     // 若元素读取数大于阈值
17     // 则 shouldSpill
18     shouldSpill = shouldSpill || _elementsRead > numElementsForceSpillThreshold
19     if (shouldSpill) {
20         // 跟新 Spill 次数
21         _spillCount += 1
22         logSpillage(currentMemory)
23         // Spill操作
24         spill(collection)
25         // 元素读取数 清零
26         _elementsRead = 0
27         // 增加Spill的内存计数
28         // 释放内存
29         _memoryBytesSpilled += currentMemory
30         releaseMemory()
31     }
32     shouldSpill
33 }

```

util.collection.Spillable.spill

将内存中的集合spill到一个有序文件中。之后SortShuffleWriter.write 中会调用 sorter.writePartitionedFile 来merge 它们

```

1  override protected[this] def spill(collection: WritablePartitionedPairCollection[K, C]): Unit = {
2  // 生成内存中集合的迭代器,
3  // 这部分我们之后会深入讲解
4  val inMemoryIterator = collection.destructiveSortedWritablePartitionedIterator(comparator)
5  // 生成spill文件,
6  // 并将其加入数组
7  val spillFile = spillMemoryIteratorToDisk(inMemoryIterator)
8  spills += spillFile
9  }

```

util.collection.ExternalSorter.spillMemoryIteratorToDisk

```

1  private[this] def spillMemoryIteratorToDisk(inMemoryIterator: WritablePartitionedIterator)
2  : SpilledFile = {
3  // 生成临时文件 及 blockId
4  val (blockId, file) = diskBlockManager.createTempShuffleBlock()
5
6  // 这些值在每次flush后会被重置
7  var objectsWritten: Long = 0
8  val spillMetrics: ShuffleWriteMetrics = new ShuffleWriteMetrics
9  val writer: DiskBlockObjectWriter =
10     blockManager.getDiskWriter(blockId, file, serInstance, fileBufferSize, spillMetrics)
11
12  // 按写入磁盘的顺序记录分支的大小
13  val batchSize = new ArrayBuffer[Long]
14
15  // 记录每个分区有多少元素
16  val elementsPerPartition = new Array[Long](numPartitions)
17
18  // Flush writer 内容到磁盘,
19  // 并更新相关变量
20  def flush(): Unit = {
21     val segment = writer.commitAndGet()
22     batchSize += segment.length
23     _diskBytesSpilled += segment.length
24     objectsWritten = 0
25  }
26
27  var success = false
28  try {
29  // 遍历内存集合
30     while (inMemoryIterator.hasNext) {
31         val partitionId = inMemoryIterator.nextPartition()
32         require(partitionId >= 0 && partitionId < numPartitions,
33             s"partition Id: ${partitionId} should be in the range [0, ${numPartitions})")
34         inMemoryIterator.writeNext(writer)
35         elementsPerPartition(partitionId) += 1
36         objectsWritten += 1
37
38         // 当写入的元素个数 到达 批量序列化尺寸,
39         // flush
40         if (objectsWritten == serializerBatchSize) {
41             flush()
42         }
43     }
44     if (objectsWritten > 0) {
45         // 遍历结束后还有写入
46         // flush
47         flush()
48     } else {
49         writer.revertPartialWritesAndClose()

```

```

50     }
51     success = true
52 } finally {
53     if (success) {
54         writer.close()
55     } else {
56         writer.revertPartialWritesAndClose()
57         if (file.exists()) {
58             if (!file.delete()) {
59                 logWarning(s"Error deleting ${file}")
60             }
61         }
62     }
63 }
64
65 SpilledFile(file, blockId, batchSize.toArray, elementsPerPartition)
66 }

```

排序

我们再次回到，`SortShuffleWriter.write` 中：

```

1    // 在外部排序中，
2    // 有部分结果可能在内存中
3    // 另外部分结果在一个或多个文件中
4    // 需要将它们merge成一个大文件
5    val partitionLengths = sorter.writePartitionedFile(blockId, tmp)

```

调用栈：

- util.collection.writePartitionedFile
 - util.collection.ExternalSorter.destructiveSortedWritablePartitionedIterator
 - util.collection.ExternalSorter.partitionedIterator
 - partitionedDestructiveSortedIterator

util.collection.ExternalSorter.writePartitionedFile

我们先来深入看下 `writePartitionedFile`，将数据加入这个 `ExternalSorter` 中，写入一个磁盘文件：

```

1    def writePartitionedFile(
2        blockId: BlockId,
3        outputFile: File): Array[Long] = {
4
5        // 跟踪输出文件的位置
6        val lengths = new Array[Long](numPartitions)
7        val writer = blockManager.getDiskWriter(blockId, outputFile, serInstance, fileBufferSize,
8            context.taskMetrics().shuffleWriteMetrics)
9
10       if (spills.isEmpty) {
11           // 当只有内存中有数据时
12           val collection = if (aggregator.isDefined) map else buffer
13           val it = collection.destructiveSortedWritablePartitionedIterator(comparator)
14           while (it.hasNext) {
15               val partitionId = it.nextPartition()
16               while (it.hasNext && it.nextPartition() == partitionId) {
17                   it.writeNext(writer)
18               }
19               val segment = writer.commitAndGet()
20               lengths(partitionId) = segment.length

```



```

21     }
22   } else {
23     // 否则必须进行merge-sort
24     // 得到一个分区迭代器
25     // 并且直接把所有数据写入
26     for ((id, elements) <- this.partitionedIterator) {
27       if (elements.hasNext) {
28         for (elem <- elements) {
29           writer.write(elem._1, elem._2)
30         }
31         val segment = writer.commitAndGet()
32         lengths(id) = segment.length
33       }
34     }
35   }
36
37   writer.close()
38   context.taskMetrics().incMemoryBytesSpilled(memoryBytesSpilled)
39   context.taskMetrics().incDiskBytesSpilled(diskBytesSpilled)
40   context.taskMetrics().incPeakExecutionMemory(peakMemoryUsedBytes)
41
42   lengths
43 }

```

util.collection.ExternalSorter.destructiveSortedWritablePartitionedIterator

在 `writePartitionedFile` 使用 `destructiveSortedWritablePartitionedIterator` 生成了迭代器:

```
1 val it = collection.destructiveSortedWritablePartitionedIterator(comparator)
```

在 [上篇博文](#) 中提到 `util.collection.Spillable.spill` 中也使用到了它:

```
1 val inMemoryIterator = collection.destructiveSortedWritablePartitionedIterator(comparator)
```

我们来看下 `destructiveSortedWritablePartitionedIterator` :

```

1 def destructiveSortedWritablePartitionedIterator(keyComparator: Option[Comparator[K]])
2   : WritablePartitionedIterator = {
3   // 生成迭代器
4   val it = partitionedDestructiveSortedIterator(keyComparator)
5   new WritablePartitionedIterator {
6     private[this] var cur = if (it.hasNext) it.next() else null
7
8     def writeNext(writer: DiskBlockObjectWriter): Unit = {
9       writer.write(cur._1._2, cur._2)
10      cur = if (it.hasNext) it.next() else null
11    }
12
13    def hasNext(): Boolean = cur != null
14
15    def nextPartition(): Int = cur._1._1
16  }
17 }

```

可以看到 `WritablePartitionedIterator` 相当于 `partitionedDestructiveSortedIterator` 所返回的迭代器的代理类。`destructiveSortedWritablePartitionedIterator` 并不返回值, 而是将 `DiskBlockObjectWriter` 传入, 再进行写。我们先把 `partitionedDestructiveSortedIterator` 放一下, 往下看。

util.collection.ExternalSorter.partitionedIterator

和另外一个分支不同, 这个分支是调用 `partitionedIterator` 得到分区迭代器, 并且直接把所有数据写入。我们来深入看看 `partitionedIterator` :

```

1  def partitionedIterator: Iterator[(Int, Iterator[Product2[K, C]])] = {
2    val usingMap = aggregator.isDefined
3    val collection: WritablePartitionedPairCollection[K, C] = if (usingMap) map else buffer
4    if (spills.isEmpty) {
5      // 当没有spills
6      // 按我们之前的流程 不会 加入这分支
7      if (!ordering.isDefined) {
8        // 若不需要对key排序
9        // 则只对Partition进行排序
10     groupByPartition(destructiveIterator(collection.partitionedDestructiveSortedIterator(None)))
11   } else {
12     // 否则需要对partition和key 进行排序
13     groupByPartition(destructiveIterator(
14       collection.partitionedDestructiveSortedIterator(Some(keyComparator))))
15   }
16 } else {
17   // 当有spills
18   // 需要 Merge spilled出来的那些临时文件 和 内存中的 数据
19   merge(spills, destructiveIterator(
20     collection.partitionedDestructiveSortedIterator(comparator)))
21 }
22 }

```

我们先来看下 `spills.isEmpty` 时候, 两种排序方式:

- 只对Partition进行排序:

`partitionedDestructiveSortedIterator` 中传入的是 `None`, 意思是不对 `key` 进行排序。对Partition进行排序是默认会在 `partitionedDestructiveSortedIterator` 中进行的。我们留在后面讲解。

```

1  groupByPartition(destructiveIterator(collection.partitionedDestructiveSortedIterator(None)))

```

Partition排序后, 根据 `Partition` 的聚合:

```

1  private def groupByPartition(data: Iterator[((Int, K), C)])
2    : Iterator[(Int, Iterator[Product2[K, C]])] =
3  {
4    val buffered = data.buffered
5    (0 until numPartitions).iterator.map(p => (p, new IteratorForPartition(p, buffered)))
6  }

```

`IteratorForPartition` 就是对单个 `partition` 的迭代器:

```

1  private[this] class IteratorForPartition(partitionId: Int, data: BufferedIterator[((Int, K), C)])
2    extends Iterator[Product2[K, C]]
3  {
4    override def hasNext: Boolean = data.hasNext && data.head._1._1 == partitionId
5
6    override def next(): Product2[K, C] = {
7      if (!hasNext) {
8        throw new NoSuchElementException
9      }
10     val elem = data.next()
11     (elem._1._2, elem._2)
12   }
13 }

```

- 对partition和key进行排序

```

1  groupByPartition(destructiveIterator(
2      collection.partitionedDestructiveSortedIterator(Some(keyComparator))))

```

`partitionedDestructiveSortedIterator` 中传入的是 `keyComparator` :

```

1  private val keyComparator: Comparator[K] = ordering.getOrElse(new Comparator[K] {
2      override def compare(a: K, b: K): Int = {
3          val h1 = if (a == null) 0 else a.hashCode()
4          val h2 = if (b == null) 0 else b.hashCode()
5          if (h1 < h2) -1 else if (h1 == h2) 0 else 1
6      }
7  })

```

先根据key的hashCode进行排序, 再调用 `groupByPartition` 对 `partition` 进行排序。

而对于有 `spills` 时, 我们使用 `comparator` :

```

1  private def comparator: Option[Comparator[K]] = {
2      // 若需要排序 或者 需要 聚合
3      if (ordering.isDefined || aggregator.isDefined) {
4          Some(keyComparator)
5      } else {
6          None
7      }
8  }

```

partitionedDestructiveSortedIterator

好了 接下来 我们就来看看 `partitionedDestructiveSortedIterator` 。 `partitionedDestructiveSortedIterator` 是 特质 `WritablePartitionedPairCollection` 中的方法。 `WritablePartitionedPairCollection` 由 `PartitionedAppendOnlyMap` 和 `PartitionedPairBuffer` 继承。在 `partitionedIterator` 中可以看到:

```

1  val usingMap = aggregator.isDefined
2  val collection: WritablePartitionedPairCollection[K, C] = if (usingMap) map else buffer

```

若需要聚合, 则使用 `PartitionedAppendOnlyMap`, 否则使用 `PartitionedPairBuffer`

util.collection.PartitionedPairBuffer.partitionedDestructiveSortedIterator

我们先来看下简单的 `PartitionedPairBuffer.partitionedDestructiveSortedIterator` :

```

1  override def partitionedDestructiveSortedIterator(keyComparator: Option[Comparator[K]])
2      : Iterator[((Int, K), V)] = {
3      val comparator = keyComparator.map(partitionKeyComparator).getOrElse(partitionComparator)
4      // 对数据进行排序
5      new Sorter(new KVArraySortDataFormat[(Int, K), AnyRef]).sort(data, 0, curSize, comparator)
6      iterator
7  }

```

我们可以看到上述:

```

1  val comparator = keyComparator.map(partitionKeyComparator).getOrElse(partitionComparator)

```

使用 `partitionKeyComparator` 将原来的 `comparator` 给替换了。 `partitionKeyComparator` 就是 `partition` 和 `key` 二次排序, 如果传入的 `keyComparator` 为 `None`, 那就是只对 `Partition` 进行排序:

```

1  def partitionKeyComparator[K](keyComparator: Comparator[K]): Comparator[(Int, K)] = {
2      new Comparator[(Int, K)] {
3          override def compare(a: (Int, K), b: (Int, K)): Int = {
4              val partitionDiff = a._1 - b._1

```

```

5         if (partitionDiff != 0) {
6             partitionDiff
7         } else {
8             keyComparator.compare(a._2, b._2)
9         }
10    }
11 }

```

之后我们使用 **Sort** 等对数据进行排序，其中用到了 **TimSort**，在以后博文中，我们会深入讲解。

最后返回迭代器 **iterator**，其实就是简单的按一对一对的去遍历数据：

```

1 private def iterator(): Iterator[((Int, K), V)] = new Iterator[((Int, K), V)] {
2     var pos = 0
3
4     override def hasNext: Boolean = pos < curSize
5
6     override def next(): ((Int, K), V) = {
7         if (!hasNext) {
8             throw new NoSuchElementException
9         }
10        val pair = (data(2 * pos).asInstanceOf[(Int, K)], data(2 * pos + 1).asInstanceOf[V])
11        pos += 1
12        pair
13    }
14 }
15 }

```

util.collection.PartitionedAppendOnlyMap.partitionedDestructiveSortedIterator

```

1 def partitionedDestructiveSortedIterator(keyComparator: Option[Comparator[K]])
2   : Iterator[((Int, K), V)] = {
3     val comparator = keyComparator.map(partitionKeyComparator).getOrElse(partitionComparator)
4     destructiveSortedIterator(comparator)
5 }

```

util.collection.PartitionedAppendOnlyMap.destructiveSortedIterator

```

1 def destructiveSortedIterator(keyComparator: Comparator[K]): Iterator[(K, V)] = {
2     destroyed = true
3     // 向左整理
4     var keyIndex, newIndex = 0
5     while (keyIndex < capacity) {
6         if (data(2 * keyIndex) != null) {
7             data(2 * newIndex) = data(2 * keyIndex)
8             data(2 * newIndex + 1) = data(2 * keyIndex + 1)
9             newIndex += 1
10        }
11        keyIndex += 1
12    }
13    assert(curSize == newIndex + (if (haveNullValue) 1 else 0))
14
15    new Sorter(new KVertexArraySortDataFormat[K, AnyRef]).sort(data, 0, newIndex, keyComparator)
16
17    // 返回新的 Iterator
18    new Iterator[(K, V)] {
19        var i = 0
20        var nullValueReady = haveNullValue
21        def hasNext: Boolean = (i < newIndex || nullValueReady)
22        def next(): (K, V) = {
23            if (nullValueReady) {
24                nullValueReady = false

```

```
25         (null.asInstanceOf[K], nullValue)
26     } else {
27         val item = (data(2 * i).asInstanceOf[K], data(2 * i + 1).asInstanceOf[V])
28         i += 1
29         item
30     }
31 }
32 }
33 }
```