

Spark MLlib StreamingKmeans实时KMeans聚类算法源代码解读

最近花了一段时间看了一下聚类算法中的StreamingKMeans算法。在Spark MLlib的聚类算法中实现多个聚类算法，KMeans也包含多个版本，其中一个是通过Spark平台来实现Kmeans，还有一个是通过Spark-streaming平台来实现Kmeans。这个类似于数据通过一个实时平台比如说Flume，或者是Kafka或者是通过Socket来发送到Spark-Streaming，然后Spark-Streaming来进行实时的数据聚类。

关于Kmeans的思想，可以参考我之前的一篇博客。

<http://blog.csdn.net/stevekangpei/article/details/73380638>

Spark-Streaming实现实时聚类的方法还是比较简单的。在StreamingKMeans的代码注释里面包含了一样一段话。

```
1  /**
2   * StreamingKMeansModel extends MLlib's KMeansModel for streaming
3   * algorithms, so it can keep track of a continuously updated weight
4   * associated with each cluster, and also update the model by
5   * doing a single iteration of the standard k-means algorithm.
6   *
7   *
8   * The update algorithm uses the "mini-batch" KMeans rule,
9   * generalized to incorporate forgetfulness (i.e. decay).
10  * The update rule (for each cluster) is:
11  *
12  * {{{
13  * c_{t+1} = [(c_t * n_t * a) + (x_t * m_t)] / [n_t + m_t]
14  * n_{t+1} = n_t * a + m_t
15  * }}}
16  *
17  * Where c_t is the previously estimated centroid for that cluster,
18  * n_t is the number of points assigned to it thus far, x_t is the centroid
19  * estimated on the current batch, and m_t is the number of points assigned
20  * to that centroid in the current batch.
21  *
22  * The decay factor 'a' scales the contribution of the clusters as estimated thus far,
23  * by applying a as a discount weighting on the current point when evaluating
```

```

24 * new incoming data. If a=1, all batches are weighted equally. If a=0, new centroids
25 * are determined entirely by recent data. Lower values correspond to
26 * more forgetting.
27 *
28 * Decay can optionally be specified by a half life and associated
29 * time unit. The time unit can either be a batch of data or a single
30 * data point. Considering data arrived at time t, the half life h is defined
31 * such that at time t + h the discount applied to the data from t is 0.5.
32 * The definition remains the same whether the time unit is given
33 * as batches or points.
34 */

```

```

1 StreamingKMeansModel继承自MLlib的KMeansModel来作为实时处理的算法。因此它可以持续的追踪着和每一个cluster
2
3 更新算法采用“mini-batch” KMeans 方法，同时也包含了消失因子(decay)。
4 更新的法则如下。
5
6 {{{
7  $c_{t+1} = [(c_t * n_t * a) + (x_t * m_t)] / [n_t + m_t]$ 
8  $n_{t+t} = n_t * a + m_t$ 
9 }}}
10
11 其中 $c_t$  是之前估计的cluster的中心点， $n_t$ 表示的是这个中心点的点的个数。 $x_t$ 是现在当前这个batch的中心点，
12 那么可以理解， $c_{t+1}$ 就是当前处理后的中心点， $n_{t+t}$ 表示的就是当前中心点的周边的点的个数。
13
14 decay因子 a 表示的是之前的数据集的下降因子。如果a = 1 的话，表示的是所有的batch的数据的权重都 是一样的
15 越多的下降因子。
16
17 decay因子可以通过一个 time unit 来指定，time unit 可以是一个 数据的batch，也可以是一个单个的数据点。如身

```

实时Kmeans聚类算法最重要的方法是实现了这个注释中的那个数学公式来进行聚类。这个算法在upc方法里面。StreamingKMeans有两个类，一个是StreamingKMeansModel，另一个是StreamingKMeans。

```

1 class StreamingKMeansModel @Since("1.2.0") (
2     @Since("1.2.0") override val clusterCenters: Array[Vector],
3     @Since("1.2.0") val clusterWeights: Array[Double])
4 extends KMeansModel(clusterCenters) with Logging {

```

```
5 //第一个参数clusterCenters表示的是数据的聚类中心点，每个元素为Vector。
6 //第二个参数指的是clusterWeights，通过对上面的算法的分析，我认为这个参数表示的是每个数据中心点这个类别
7
```

接下来是一个update方法，实时聚类算法也是在这个方法中实现的。

```
1 def update(data: RDD[Vector], decayFactor: Double, timeUnit: String): StreamingKMeansModel = {
2
3     // find nearest cluster to each point
4     val closest = data.map(point => (this.predict(point), (point, 1L)))
5     //这个表示的是通过当前的model找到现在这个batch数据集中地这个点所在的聚类中心的index，
6     //同时返回一个元祖格式的数据(index: Int, (point: Vector, 1: Long))。
7     //这个表示的是当前这个batch数据里面的每个数据所属的类别，和这个点出现的次数，
8     //注意我们之前的那个注释，这个算法中需要用到每个中心点的周围的点的个数。
9
10    // get sums and counts for updating each cluster
11    //这里定义了一个函数，这个函数对两个(Vector, Long)类型的元祖进行一个合并操作，
12    //返回的类型也是一个元祖的类型。
13    val mergeContribs: ((Vector, Long), (Vector, Long)) => (Vector, Long) = (p1, p2) => {
14        BLAS.axpy(1.0, p2._1, p1._1)
15        //这个做一个向量的加法操作，
16        //表示p1._1 = p1._1 + 1.0 * p2._1
17        (p1._1, p1._2 + p2._2) //然后返回的是(两个向量的和，两个向量出现的次数之和)。
18    }
19    val dim = clusterCenters(0).size //这个表示的是聚类中心点的维度的数目。
20
21    val pointStats: Array[(Int, (Vector, Long))] = closest
22        .aggregateByKey((Vectors.zeros(dim), 0L))(mergeContribs, mergeContribs)
23        .collect()
24
25    //这个里面用到了一个aggregateByKey算子，这个算子在我之前的博客中讲过，
26    //第一个参数表示的是zeroValue，初始值，这个会作用到第一个函数但是不会作用到第二个函数。
27    //第一个函数值得是对每个分区进行merge操作，就是我们之前定义的函数，
28    //第二个表示的是对每个分区作用后的结果进行操作。
29    //最后返回的结果通过collect返回一个数组，每个元素的类型为(Int (表示这个所有batch的数据出现在第int个
30    //(表示的是出现在这个中心点附近的所有的向量的矢量和，出现的向量的次数))
31
32    //这个表示的是用来判断discount 的类型，是batch类型还是points类型，如果是batch类型，则直接使用decayF
33    val discount = timeUnit match {
```

```

34     case StreamingKMeans.BATCHES => decayFactor
35     case StreamingKMeans.POINTS =>
36         val numNewPoints = pointStats.view.map { case (_, (_, n)) =>
37             n
38         }.sum
39         math.pow(decayFactor, numNewPoints)
40     }
41
42     // apply discount to weights
43     BLAS.scal(discount, Vectors.dense(clusterWeights))
44     // 这个计算好的discount的值被用来scale之前cluster的权重，这个和之前的第一个公式的第一部分是类似的。
45     //  $c_{t+1} = [(c_t * n_t * a) + (x_t * m_t)] / [n_t + m_t]$ 
46     //即将decay值作用于n_t(之前点出现的次数)。
47
48     // implement update rule
49     //然后是实现这个update算法
50     pointStats.foreach { case (label, (sum, count)) =>
51         val centroid = clusterCenters(label) //获取对应的中心点，
52         val updatedWeight = clusterWeights(label) + count
53         //更新现在的权重值，其实指的就是n_t + m_t
54
55         val lambda = count / math.max(updatedWeight, 1e-16)
56         //将count/updatedWeight作为现在的lamda值。这个表示的就是
57
58         clusterWeights(label) = updatedWeight //同时更新现在的clusterWeights数组。
59         BLAS.scal(1.0 - lambda, centroid)
60         BLAS.axpy(lambda / count, sum, centroid)
61         //这个指的就是对所有的点的向量和除以总的出现的次数来作为新的centroid。
62
63         // display the updated cluster centers
64         //接下来打印出聚类中心点。
65         val display = clusterCenters(label).size match {
66             case x if x > 100 => centroid.toArray.take(100).mkString("[", ",", "...")
67             case _ => centroid.toArray.mkString("[", ",", "]")
68         }
69
70         logInfo(s"Cluster $label updated with weight $updatedWeight and centroid: $display")
71     }
72
73     // Check whether the smallest cluster is dying. If so, split the largest cluster.
74     //接下来判断有没有哪个点正在消失，如果消失的话，那么歼最大的cluster进行分开。
75     val weightsWithIndex = clusterWeights.view.zipWithIndex

```

```

76 //这个表示将每个数据集和它出现的索引zip起来。然后找到最大的和最小的数据聚类。
77 val (maxWeight, largest) = weightsWithIndex.maxBy(_._1)
78
79 val (minWeight, smallest) = weightsWithIndex.minBy(_._1)
80 if (minWeight < 1e-8 * maxWeight) {
81 //如果最小的太小的话
82     logInfo(s"Cluster $smallest is dying. Split the largest cluster $largest into two.")
83     val weight = (maxWeight + minWeight) / 2.0
84     //将大数据聚类进行拆分，然后重新分给小数据聚类 and 大数据聚类。
85     clusterWeights(largest) = weight
86     clusterWeights(smallest) = weight
87     val largestClusterCenter = clusterCenters(largest)
88     val smallestClusterCenter = clusterCenters(smallest)
89     var j = 0
90     while (j < dim) {
91         val x = largestClusterCenter(j)
92         val p = 1e-14 * math.max(math.abs(x), 1.0)
93         largestClusterCenter.toBreeze(j) = x + p
94         smallestClusterCenter.toBreeze(j) = x - p
95         j += 1
96     }
97 }
98
99 this //最后返回这个新的模型。
100 }
101 }
102

```

这个基本的类包含的是StreamingKmeansModel的初始化。

k表示的是中心点的个数，decayFactor表示的是下降因子，timeUnit有两种类型为batch和points。

```

1 class StreamingKMeans @Since("1.2.0") (
2     @Since("1.2.0") var k: Int,
3     @Since("1.2.0") var decayFactor: Double,
4     @Since("1.2.0") var timeUnit: String) extends Logging with Serializable {

```

然后也包含一些predict方法。

```

@Since("1.4.0")
def predictOn(data: JavaDStream[Vector]): JavaDStream[java.lang.Integer] = {
  JavaDStream.fromDStream(predictOn(data.dstream).asInstanceOf[DStream[java.lang.Integer]])
}

/**
 * Use the model to make predictions on the values of a DStream and carry over its keys.
 *
 * @param data DStream containing (key, feature vector) pairs
 * @tparam K key type
 * @return DStream containing the input keys and the predictions as values
 */
@Since("1.2.0")
def predictOnValues[K: ClassTag](data: DStream[(K, Vector)]): DStream[(K, Int)] = {
  assertInitialized()
  data.mapValues(model.predict)
}

/**
 * Java-friendly version of `predictOnValues`.
 */
@Since("1.4.0")
def predictOnValues[K] (
  data: JavaPairDStream[K, Vector]): JavaPairDStream[K, java.lang.Integer] = {
  implicit val tag = fakeClassTag[K]
  JavaPairDStream.fromPairDStream(
    predictOnValues(data.dstream).asInstanceOf[DStream[(K, java.lang.Integer)])
  )
}

```