

龙猫先生

5-Spark高级数据分析-第五章 基于K均值聚类的网络流量异常检测

据我们所知，有‘已知的已知’，有些事，我们知道我们知道；我们也知道，有‘已知的未知’，也就是说，有些事，我们现在知道我们不知道。但是，同样存在‘未知的不知’——有些事，我们不知道我们不知道。

上一章中分类和回归都属于监督学习。当目标值是未知时，需要使用非监督学习，非监督学习不会学习如何预测目标值。但是，它可以学习数据的结构并找出相似输入的群组，或者学习哪些输入类型可能出现，哪些类型不可能出现。

5.1 异常检测

异常检测常用于检测欺诈、网络攻击、服务器及传感设备故障。在这些应用中，我们要能够找出以前从未见过的新型异常，如新欺诈方式、新入侵方法或新服务器故障模式。

5.2 K均值聚类

聚类是最有名的非监督学习算法，K均值聚类是应用最广泛的聚类算法。它试图在数据集中找出k个簇群。在K均值算法中数据点相互距离一般采用欧氏距离。

在K均值算法中簇群其实是一个点，即组成该簇的所有点的中心。数据点其实就是由所有数值型特征组成的特征向量，简称向量。

簇群的中心称为质心，它是簇群中所有点的算术平均值，因此算法取名K均值。算法开始时选择一些数据点作为簇群的质心。然后把每个数据点分配给最近的质心。接着对每个簇计算该簇所有数据点的平均值，并将其作为该簇的新质心。然后不断重复这个过程。

5.3 网络入侵

统计对各个端口在短时间内被远程访问的次数，就可以得到一个特征，该特征可以很好地预测端口扫描攻击。检测网络入侵是要找到与以往见过的连接不通的连接。K均值可根据每个网络连接的统计属性进行聚类，结果簇定义了历史连接类型，帮我们界定了正常的连接的区域。任何在区域之外的点都是不正常的。

5.4 KDD Cup 1999数据集

KDD Cup是数据挖掘竞赛，由ACM特别兴趣小组举办。1999年主题为网络入侵。

数据下载地址：<http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

百度云：<http://pan.baidu.com/s/1cFqnRS>

数据集大小为108，每个连接信息包括发送的字节数、登录次数、TCP错误数等。数据集为CSV格式，每个连接占一行，包括38个特征。

我们关心的问题是找到“未知”的攻击。

5.5 初步尝试聚类

加载数据并查看有哪些类别标号及每类样本有多少：

Scala：

```
1 val rawData = sc.textFile("D:/Workspace/AnalysisWithSpark/src/main/java/advanced/chapter5/kddcu
2 rawData.map(_._split(',').last).countByValue().toSeq.sortBy(_._2).reverse.foreach(println)
```

公告

昵称：[龙猫先生](#)
园龄：[2年8个月](#)
粉丝：[15](#)
关注：[3](#)
[+加关注](#)

导航

[博客园](#)
[首页](#)
[新随笔](#)
[联系](#)
[订阅](#) [XML](#)
[管理](#)

< 2018年1月 >						
日	一	二	三	四	五	六
31	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	1	2	3
4	5	6	7	8	9	10

统计

随笔 - 168
文章 - 0
评论 - 5
引用 - 0

搜索

常用链接

[我的随笔](#)
[我的评论](#)
[我的参与](#)
[最新评论](#)
[我的标签](#)

我的标签

[Spark](#)(5)
[Mysql](#)(1)

随笔档案

[2017年2月](#) (1)
[2016年8月](#) (7)
[2015年11月](#) (10)
[2015年8月](#) (5)
[2015年6月](#) (1)
[2015年2月](#) (1)
[2015年1月](#) (1)
[2014年9月](#) (1)
[2014年8月](#) (5)
[2014年7月](#) (1)
[2013年11月](#) (1)

Java :

```

1 //初始化SparkConf
2 SparkConf sc = new
SparkConf().setMaster("local").setAppName("AnomalyDetectionInNetworkTraffic");
3 System.setProperty("hadoop.home.dir", "D:/Tools/hadoop-2.6.4");
4 JavaSparkContext jsc = new JavaSparkContext(sc);
5
6 //读入数据
7 JavaRDD<String> rawData
=jsc.textFile("src/main/java/advanced/chapter5/kddcup.data/kddcup.data.corrected");
8
9 //查看有哪些类别标号及每类样本有多少
10 ArrayList<Entry<String, Long>> lineList = new ArrayList<>(rawData.map(line ->
line.split(",") [line.split(",").length-1]).countByValue().entrySet());
11 Collections.sort(lineList, (m1, m2) -> m2.getValue().intValue() -
m1.getValue().intValue());
12 lineList.forEach(line -> System.out.println(line.getKey() + "," + line.getValue()));

```

结果 :

```

smurf.,2807886
neptune.,1072017
normal.,972781
satan.,15892
ipsweep.,12481
portsweep.,10413
nmap.,2316
back.,2203
warezclient.,1020
teardrop.,979
pod.,264
guess_passwd.,53
buffer_overflow.,30
land.,21
warezmaster.,20
imap.,12
rootkit.,10
loadmodule.,9
ftp_write.,8
multihop.,7
phf.,4
perl.,3
spy.,2

```

看来用Scala一行能写完的代码用Java还是比较麻烦的。

下面将CSV格式的行拆成列，删除下标从1开始的三个类别型列和最后的标号列。

Scala :

```

1 import org.apache.spark.mllib.linalg._
2 val labelsAndData = rawData.map { line =>
3     val buffer = line.split(',').toBuffer

```

2013年10月 (1)
 2013年8月 (1)
 2013年7月 (1)
 2013年6月 (1)
 2013年5月 (4)
 2013年4月 (5)
 2013年3月 (4)
 2013年2月 (3)
 2013年1月 (9)
 2012年12月 (23)
 2012年11月 (3)
 2012年10月 (13)
 2012年9月 (16)
 2012年7月 (20)
 2012年6月 (12)
 2012年5月 (18)

最新评论

1. Re:2-Spark高级数据分析-第二章 用Scala和Spark进行数据分析
 不错，正在学习中

--loulley

2. Re:2-Spark高级数据分析-第二章 用Scala和Spark进行数据分析

基于Python Spark的大数据分析

课程观看地址：

xuetuwuyoucom/course/173

课程来自学途无忧网：

www.xuetuwuyoucom

--学习童

3.
 Re:VS2008 f:/dd/vctools/vc7libs/s
 误调试

@龙猫先生好像没用啊，我只是在弹出的对话框中 加了一个黑色显示视频区域。...

--Riven^

4.
 Re:VS2008 f:/dd/vctools/vc7libs/s
 误调试

@Riven^ 在出问题的控件中增加一个
 DrawItem函数Overwrite：重写...

--龙猫先生

5.
 Re:VS2008 f:/dd/vctools/vc7libs/s
 误调试

大哥 这怎么重写啊

--Riven^

阅读排行榜

1. 2-Spark高级数据分析-第二章 用Scala和Spark进行数据分析(4971)

2. Sigar使用(3518)

3. Mysql使用mysqldump按时间导出时的一个注意事项(2174)

4. 4-Spark高级数据分析-第四章 用决策树算法预测森林植被(1891)

5. 3-Spark高级数据分析-第三章 音乐推荐和Audioscrobbler数据集(1331)

6. 5-Spark高级数据分析-第五章 基于k均值聚类的网络流量异常检测(1069)

7. 使用SQL语法来查询

Elasticsearch: Elasticsearch-SQL插件(1004)

8. 1-Spark高级数据分析-第一章 大数据分析(368)

9. 利用powerdesigner反向数据库结构，生成ER图(299)

10. 0-Spark高级数据分析-读书笔记(270)

```

4     buffer.remove(1, 3)
5     val label = buffer.remove(buffer.length-1)
6     val vector = Vectors.dense(buffer.map(_.toDouble).toArray)
7     (label,vector)
8 }
9 val data = labelsAndData.values.cache()

```

评论排行榜

1. VS2008 f:/dd/vctools/vc7libs/ship 误调试(3)
2. 2-Spark高级数据分析-第二章 用Scala和Spark进行数据分析(2)

推荐排行榜

1. 4-Spark高级数据分析-第四章 用决策树算法预测森林植被(2)
2. 3-Spark高级数据分析-第三章 音乐推荐和Audioscrobbler数据集(2)
3. 2-Spark高级数据分析-第二章 用Scala和Spark进行数据分析(2)
4. 0-Spark高级数据分析-读书笔记(1)
5. Sigar使用(1)

Java :



```

1 //删除下标从1开始的三个类别列和最后的标号列
2 JavaRDD<Tuple2<String, Vector>> labelsAndData = rawData.map(line -> {
3     String[] lineArrya = line.split(",");
4     double[] vectorDouble = new double[lineArrya.length-4];
5     for (int i = 0, j=0; i < lineArrya.length; i++) {
6         if(i==1 || i==2 || i==3 || i==lineArrya.length-1) {
7             continue;
8         }
9         vectorDouble[j] = Double.parseDouble(lineArrya[i]);
10        j++;
11    }
12    String label = lineArrya[lineArrya.length-1];
13    Vector vector = Vectors.dense(vectorDouble);
14    return new Tuple2<String, Vector>(label,vector);
15 });
16
17 RDD<Vector> data = JavaRDD.toRDD(labelsAndData.map(f -> f._2));

```



对数据进行聚类

Scala :

```

1 import org.apache.spark.mllib.clustering._
2 val kmeans = new KMeans()
3 val model = kmeans.run(data)
4 model.clusterCenters.foreach(println)

```

Java :



```

1 //聚类
2 KMeans kmeans = new KMeans();
3 KMeansModel model = kmeans.run(data);
4
5 //聚类结果
6 Arrays.asList(model.clusterCenters()).forEach(v -> System.out.println(v.toJson()));

```



结果 :

```

{"type":1,"values":
[48.34019491959669,1834.6215497618625,826.2031900016945,5.7161172049003456
E-6,6.487793027561892E-4,7.961734678254053E-

```

```
6,0.012437658596734055,3.205108575604837E-
5,0.14352904910348827,0.00808830584493399,6.818511237273984E-
5,3.6746467745787934E-
5,0.012934960793560386,0.0011887482315762398,7.430952366370449E-
5,0.0010211435092468404,0.0,4.082940860643104E-7,8.351655530445469E-
4,334.9735084506668,295.26714620807076,0.17797031701994304,0.1780369894027
2675,0.05766489875327384,0.05772990937912762,0.7898841322627527,0.02117961
0609915762,0.02826081009629794,232.98107822302248,189.21428335201279,0.753
713389800417,0.030710978823818437,0.6050519309247937,0.006464107887632785,
0.1780911843182427,0.17788589813471198,0.05792761150001037,0.0576592214240
0437]]}

{"type":1,"values":
[10999.0,0.0,1.309937401E9,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,
,0.0,0.0,0.0,0.0,1.0,1.0,0.0,0.0,1.0,1.0,1.0,0.0,0.0,255.0,1.0,0.0,0.65,1.
0,0.0,0.0,0.0,1.0,1.0]]}
```

程序输出两个向量，代表K均值将数据聚类成k=2个簇。对本章的数据集，我们知道连接的类型有23个，因此程序肯定没能准确刻画出数据中的不同群组。

查看两个簇中分别包含哪些类型的样本。

Scala :

```
1 val clusterLabelCount = labelsAndData.map { case (label,datum) =>
2     val cluster = model.predict(datum)
3     (cluster,label)
4 }.countByValue
5 clusterLabelCount.toSeq.sorted.foreach {
6     case ((cluster,label),count) =>
7         println(f"$cluster%1s$label%18s$count%8s")
8 }
```

Java :



```
1 ArrayList<Entry<Tuple2<Integer, String>, Long>> clusterLabelCount = new
ArrayList<Entry<Tuple2<Integer, String>, Long>>(labelsAndData.map( v -> {
2     int cluster = model.predict(v._2);
3     return new Tuple2<Integer, String>(cluster, v._1);
4 }).countByValue().entrySet());
5
6 Collections.sort(clusterLabelCount, (m1, m2) -> m2.getKey()._1-m1.getKey()._1);
7 clusterLabelCount.forEach(t -> System.out.println(t.getKey()._1 + "\t" + t.getKey()._2
+ "\t\t" + t.getValue()));
```



结果 :

```
1 portsweep. 1
0 portsweep. 10412
0 rootkit. 10
0 buffer_overflow. 30
0 phf. 4
0 pod. 264
```

```

0 perl. 3
0 spy. 2
0 ftp_write. 8
0 nmap. 2316
0 ipsweep. 12481
0 imap. 12
0 warezmaster. 20
0 satan. 15892
0 teardrop. 979
0 smurf. 2807886
0 neptune. 1072017
0 loadmodule. 9
0 guess_passwd. 53
0 normal. 972781
0 land. 21
0 multihop. 7
0 warezclient. 1020
0 back. 2203

```

结果显示聚类根本没有任何作用。簇1只有一个数据点！

5.6 K的选择

计算两点距离函数：

Scala：

```

1 def distance(a: Vector, b: Vector) =
2   math.sqrt(a.toArray.zip(b.toArray).
3     map(p => p._1 - p._2).map(d => d * d).sum)

```

Java：



```

1 public static double distance(Vector a, Vector b){
2   double[] aArray = a.toArray();
3   double[] bArray = b.toArray();
4   ArrayList<Tuple2<Double, Double>> ab = new ArrayList<Tuple2<Double, Double>>();
5   for (int i = 0; i < a.toArray().length; i++) {
6     ab.add(new Tuple2<Double, Double>(aArray[i],bArray[i]));
7   }
8   return Math.sqrt(ab.stream().map(x -> x._1-x._2).map(d -> d*d).reduce((r,e) -> r=
r+e).get());
9 }

```



计算数据点到簇质心距离函数：

Scala：

```

1 def distToCentroid(datum: Vector, model: KMeansModel) = {
2   val cluster = model.predict(datum)
3   val centroid = model.clusterCenters(cluster)
4   distance(centroid, datum)
5 }

```

Java :

```
1 public static double distToCentroid(Vector datum, KMeansModel model) {
2     int cluster = model.predict(datum);
3     Vector[] centroid = model.clusterCenters();
4     return distance(centroid[cluster], datum);
5 }
```

给定k值的模型的平均质心距离函数：

Scala :

```
1 import org.apache.spark.rdd._
2 def clusteringScore(data: RDD[Vector], k: Int) = {
3     val kmeans = new KMeans()
4     kmeans.setK(k)
5     val model = kmeans.run(data)
6     data.map(datum => distToCentroid(datum, model)).mean()
7 }
```

Java :

```
1 public static double clusteringScore(JavaRDD<Vector> data, int k) {
2     KMeans kmeans = new KMeans();
3     kmeans.setK(k);
4     KMeansModel model = kmeans.run(JavaRDD.toRDD(data));
5     return data.mapToDouble(datum -> distToCentroid(datum, model)).stats().mean();
6 }
```

对K从5到40进行评估：

Scala :

```
1 (5 to 40 by 5).map(k => (k, clusteringScore(data, k))).foreach(println)
```

Java :

```
1 List<Double> list = Arrays.asList(new Integer[]{1, 2, 3, 4, 5, 6, 7, 8}).stream().map(k
-> clusteringScore(labelsAndData.map(f -> f._2), k*5)).collect(Collectors.toList()); 2 3
list.forEach(System.out::println);
```

要算很久，结果：

```
1938.8583418059206
1686.4806829850777
1440.0646239087368
1305.763038353858
964.3070891182899
878.7358671386651
571.8923560384558
745.7857049862099
```

5.11 聚类实战

偷懒了，中间的那些和R相关还有标准化的没有写。

取k=150, 聚类结果如下：

```
149 normal. 4
148 warezclient. 590
148 guess_passwd. 52
148 nmap. 1472
148 portsweep. 378
148 imap. 9
148 ftp_write. 2
....
97 warezclient. 275
96 normal. 3
95 normal. 1
94 normal. 126
93 normal. 47
92 normal. 52196
92 loadmodule. 1
92 satan. 1
92 buffer_overflow.3
92 guess_passwd. 1
91 normal. 1
90 normal. 3
89 normal. 6
88 normal. 12388
....
16 normal. 1
15 normal. 11
14 normal. 68
13 normal. 232
12 normal. 1
11 portsweep. 1
10 portsweep. 1
9 warezclient. 59
9 normal. 1
8 normal. 1
7 normal. 1
6 portsweep. 1
5 portsweep. 1
4 portsweep. 1
3 portsweep. 2
2 portsweep. 1
1 portsweep. 1
0 smurf. 527579
0 normal. 345
```

作为示例, 我们在原始数据上进行异常检查：

Scala：

```
1 val model = ...
2 val originalAndData = ...
3 val anomalies = originalAndData.filter { case (original, datum) =>
4     val normalized = normalizeFunction(datum)
5     distToCentroid(normalized, model) > threshold
```



```
6 | }.keys
```

Java :


```

1      KMeans kmeansF = new KMeans();
2      kmeansF.setK(150);
3      KMeansModel modelF = kmeansF.run(data);
4
5      System.out.println("json:-----");
6      Arrays.asList(modelF.clusterCenters()).forEach(v ->
System.out.println(v.toJson()));
7
8      ArrayList<Entry<Tuple2<Integer, String>, Long>> clusterLabelCountF = new
ArrayList<Entry<Tuple2<Integer, String>, Long>>(labelsAndData.map( v -> {
9          int cluster = modelF.predict(v._2);
10         return new Tuple2<Integer, String>(cluster, v._1);
11     })).countByValue().entrySet());
12
13     Collections.sort(clusterLabelCountF, (m1, m2) -> m2.getKey()._1-
m1.getKey()._1);
14     clusterLabelCountF.forEach(t -> System.out.println(t.getKey()._1 +"\t"+
t.getKey()._2 +"\t\t"+ t.getValue()));
15
16     //距离中心最远的第100个点的距离
17     JavaDoubleRDD distances = labelsAndData.map(f -> f._2).mapToDouble(datum ->
distToCentroid(datum, modelF));
18     Double threshold = distances.top(100).get(99);
19
20     JavaRDD<Tuple2<String, Vector>> result = labelsAndData.filter(t ->
distToCentroid(t._2, modelF) > threshold);
21     System.out.println("result:-----");
22     result.foreach(f -> System.out.println(f._2));

```



结果如下：

[illegible]


```
[60.0,854.0,1519233.0,0.0,0.0,0.0,0.0,0.0,1.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,1.0,1.0,0.0,0.0,0.0,0.0,1.0,0.0,0.0,113.0,34.0,0.3,0.04,0.01,0.0,0.0,0.0,0.0,0.0]
[107.0,585.0,2661605.0,0.0,0.0,0.0,0.0,0.0,1.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,1.0,1.0,0.0,0.0,0.0,0.0,1.0,0.0,0.0,171.0,47.0,0.27,0.02,0.01,0.0,0.0,0.0,0.0,0.0]
```

.....

.....

5.12 小结

可以改成StreamingKmeans，它会根据增量对簇进行更新。官方文档中也只有用Scala写的代码，如果需要找Java的话，可以参考我的另外一个项目中的代码：

<https://github.com/jiangpz/LearnSpark/blob/master/src/main/java/mllib/StreamingKmeansExample.java>

标签: [Spark](#)

[好文要顶](#)[关注我](#)[收藏该文](#)

龙猫先生

关注 - 3

粉丝 - 15

[+加关注](#)

« 上一篇: [4-Spark高级数据分析-第四章 用决策树算法预测森林植被](#)

» 下一篇: [Mysql使用mysqldump按时间导出时的一个注意事项](#)

posted on 2016-08-24 17:14 [龙猫先生](#) 阅读(1069) 评论(0) [编辑](#) [收藏](#)

[刷新评论](#) [刷新页面](#) [返回顶部](#)

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问网站首页](#)。

【推荐】50万行VC++源码：大型组态工控、电力仿真CAD与GIS源码库

【推荐】加入腾讯云自媒体扶持计划，免费领取域名&服务器

【福利】限时领取，H3 BPM给你发年终奖



最新IT新闻:

- [明明开启了 支付宝收钱码为啥还没有语音提醒？](#)
- [特斯拉发布“测试版”自动感应雨刷 未来将推广到全车型](#)
- [外媒揭秘苹果、谷歌和Facebook如何跟踪用户位置](#)
- [微软再发对比视频：Edge浏览器续航比Firefox多63%](#)
- [不甘给小米代工 富士康自有品牌手机进军印度](#)
- » [更多新闻...](#)



最新知识库文章:

- [步入云计算](#)
- [以操作系统的角度述说线程与进程](#)
- [软件测试转型之路](#)
- [门内门外看招聘](#)
- [大道至简，职场上做人做事做管理](#)

» [更多知识库文章...](#)

Powered by: [博客园](#) Copyright © 龙猫先生