

小黑

English:Talk is cheap,show me the code. 中文：P话少说，放码过来。

目录视图

摘要视图

RSS 订阅

个人资料



heybiiiii

关注

发私信

访问：385506次

积分：5319

等级：BLOG > 5

排名：第5829名

原创：164篇

转载：2篇

译文：0篇

评论：111条

文章搜索

文章分类

- Hadoop (43)
- Datamining (6)
- Spark (19)
- Docker (13)
- .NET (34)
- JS (4)
- 设计模式 (1)
- Others (11)
- Linux (5)
- Unity3D (1)
- 软件测试 (10)
- HBase (11)
- Kafka (1)
- Scala (2)
- Storm (3)
- 概率论与数理统计 (0)
- 线性代数 (0)
- Drools (2)

文章存档

- 2016年09月 (1)
- 2016年07月 (6)

图灵赠书——程序员11月书单 【思考】Python这么厉害的原因竟然是！ 感恩节赠书：《深度学习》等异步社区优秀图书和作译者评选启动！ 每周荐书：京东架构、Linux内核、Python全栈

[置顶] KMeans算法检测网络异常入侵

标签：算法 KMeans Spark SparkMllib

2016-05-09 17:09

4303人阅读

评论(1)

分类：

Spark (18)

版权声明：本文为博主原创文章，未经博主允许不得转载。

目录(?)

[+]

非监督学习技术

在决策树算法预测森林植被中

我们可以体会到属于监督学习的分类和回归技术的强大，可以预测“即将发生”的事情使用监督学习技术有一个很关键的前提：需要大量的数据对模型进行训练，模型能够从已知的数据中学习规律进而预测未知的数据

然而在某些场景下，并不是都能提供监督学习所需要的样本数据来训练模型，有可能只能给出部分正确的输出，甚至一个输出都没有这种情况下，监督学习的技术就不能够使用了

此时，对应监督学习，另一种非监督学习技术就可以排上用场了

异常检查

顾名思义，异常检测就是要找出不同寻常的情况，异常是一种未知的情况，也就是说，无论何时何地，我们都无法归纳总结出所有的异常分类如果可以，那么使用监督学习技术可以轻易的将网站的每个访问划分为“正常”或者“异常”举个例子来说，我们永远不知道黑客有什么新的技术手段可以入侵你的网站系统，即使今天你有所有已知的黑客手段，但是谁知道明天又会有新的漏洞被黑客利用？所以说，当一个访问请求被处理的时候，如果使用监督学习技术，恰好这是一个异常访问的请求，又恰好这是一种全新的异常类别此时监督学习技术就会束手无策

在这种场景下，使用非监督学习技术可以有效的解决这个问题，通过学习，它们能够知道什么是正常的输入从而能够判别出新数据和历史数据的差别，注意，这里的有差别并不意味着该数据就是异常数据，只是说它和历史的正常数据有差异，值得进一步调查

所以，非监督学习并不是要将数据精确地划分到哪个类别中，而是对历史数据进行对比分析找出差异

2016年05月 (10)

2016年03月 (7)

2016年02月 (3)

展开

阅读排行

Kubernetes用户指南（一）--... (23766)

Hadoop/Spark相关面试问题... (15076)

Kubernetes用户指南（四）--... (11828)

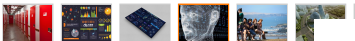
基于Go的机器学习框架... (40405)



可视化数据



人脸识别算法



推荐文章

* 【2017年11月27日】CSDN博客更新周报

* 【CSDN】邀请您来GitChat赚钱啦！

* 【GitChat】精选——JavaScript进阶指南

* 改做人工智能之前，90%的人都没能给自己定位

* TensorFlow 人脸识别网络与对抗网络搭建

* Vue 移动端项目生产环境优化

* 面试必考的计算机网络知识点梳理

最新评论

Kubernetes用户指南（一）--快速开始、...
专业大数据 : 交流学习java架构大数据，加群460570824

CDH5.3配置Kerberos+LDAP+Sentry记录
x6696 : 非常有用 解决了几个问题 多泄

搭建Drools开发环境
木子-轩 : @qq_15951857:<?xml version="1.0" encoding="UTF-8" ?>

MapReduce中的分布式缓存使用
大数据中的大叔 : 代码部分，测试过了，伪分布环境下测试不通过，完全分布环境下测试通过。我们传进去的Distribute...

在Docker中从头部署自己的Spark集群
2551 : docker commit {containerId}#会返回一个iddocker tag {id}...

asp.net权限控制的方式
woshisunzewe : http://122.112.248.118:8082/Account/Loginhttp://12...

KMeans均值聚类

聚类是最有名的非监督学习技术，它试图找到数据中的自然群组

一群特征相似而又与其他数据不同的数据点往往代表某种意义，从而将这些数据点划分为一个族群

聚类算法就是要将所有数据中的相似数据划分到同一个族群中

在网络的异常检测中，聚类算法是十分合适的，通过其将所有访问请求划分为一个个族群，将正常和异常的访问隔离开

我们可以进一步在异常的族群中分析这些数据是否属于网络入侵

K均值聚类是运行的最广泛的聚类算法，根据人为定义的一个k值，该算法会将数据聚类为k个族群

关于k值的确定需要结合业务场景和数据特性，并在反复的实验中得到一个最优

KMeans聚类的详细说明可以参考：

[mahout运行测试与数据挖掘算法之聚类分析（一）kmeans算法解析](#)

程序开发

数据集

案例中使用的是KDD Cup1999的数据集，可以在这里[下载](#)

进入下载页面后可以看到有很多数据集，本篇只用到kddcup.data.*数据文件，一个是完整数据集，解压缩之后有743M，一个是10%的数据集，解压缩之后只有45M

可以先使用10%的数据集进行代码测试，得到一个比较好的结果之后再运用到全部数据集中

数据集中的每一行代表一个网络请求，描述了该请求的所有信息，在[决策树算法预测森林植被](#)中

我们知道特征分为数值型和类别型，该数据集中也包含了这两种特征

每行数据的最后一个特征是目标特征，例如大部分请求被标记为normal表示正常访问

还有其他各种异常标记

正如之前讨论的，我们完全可以使用特征向量+目标特征的形式使用监督学习技术来训练模型进行预测

但是如果出现一个异常请求不在所有的目标类别中，这不糟糕了~

所以为了找出“未知的攻击”，我们先不在算法中使用这些目标特征

聚类的初步尝试

将要使用的数据集上传到HDFS之后，我们先来看看这些数据的基本信息：

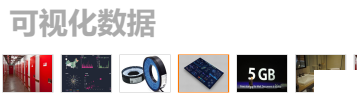
```
1 val conf = new SparkConf().setAppName("KMeans")
2 val sc = new SparkContext(conf)
3 //读取数据
4 val rawData = sc.textFile("/spark_data/ch05/kddcup.data")
5 //根据类别查看统计信息,各个类别下有多少数据
6 //根据", "分割,只保留最后的类别
7 val catStatsData = rawData.map(_.split(",").last)
8 //对类别的数目进行统计,并根据统计的数量从小打到排序
9 .countByValue().toSeq.sortBy(_._2)
10 //转换为从大到小排序
11 .reverse
12 catStatsData.foreach(println)
```

运行结果：

MapReduce高级特性

大数据中的大叔：楼主，你好！关于hadoop的全排序方面，Mapper只能控制局部排序，我们能保证每个分区中的值的顺...

多个Mapper和Reducer的Job
xun-ming : 感谢分享



总共有23种不同类型，其中smurf和neptune的攻击类型最多，竟然比normal的正常访问还要多

```
1  val labelsAndData = rawData.map { line =>
2      //buffer是一个可变列表
3      val buffer = line.split(",").toBuffer
4      //下标1-3的元素
5      buffer.remove(1, 3)
6      //最后一个元素为label
7      val label = buffer.remove(buffer.length - 1)
8      //转换为Vector
9      val vector = Vectors.dense(buffer.map(_.toDouble).toArray)
10     (label, vector)
11 }
12 //数据只用到values部分
13 val data=labelsAndData.values.cache()
```

```
1 //训练模型
2 val kmeans = new KMeans()
3 val model = kmeans.run(data)
4 //输出所有聚类中心
5 model.clusterCenters.foreach(println)
```

[illegible]



可视化数据



人脸识别算法



现在我们用这个模型来为每个数据分配族群，并输出每个聚类中心有哪些类别，各有多少个数据：

```
1 //输出每个聚类中心有哪些类别各有多少个数据
2 val clusterLabelCount = labelsAndData.map { case (label, datum) =>
3     //为样本数据划分聚类中心
4     val cluster = model.predict(datum)
5     //返回数据的中心和类别二元组
6     (cluster, label)
7 }.countByValue()
8 //排序之后格式化输出
9 clusterLabelCount.toSeq.sorted.foreach { case ((cluster, label), count) =>
10     println(f"$cluster%1s$label%1s$count")
11 }
```

由于我们没有人为设置k的值，程序就自作主张设置了k=2，但是我们知道，数个，因此这个模型肯定是错误的
可以看到输出的结果为：

3

```
1 0 back. 2203
2 0 buffer_overflow. 30
3 0 ftp_write. 8
4 0 guess_passwd. 53
5 0 imap. 12
6 0 ipsweep. 1247
7 0 land. 21
8 0 loadmodule. 9
9 0 multihop. 7
10 0 neptune. 107201
11 0 nmap. 231
12 0 normal. 97278
13 0 perl. 3
14 0 phf. 4
15 0 pod. 264
16 0 portsweep. 1039
17 0 rootkit. 10
18 0 satan. 1589
19 0 smurf. 280790
20 0 spy. 2
21 0 teardrop. 979
22 0 warezclient. 1020
23 0 warezmaster. 20
24 1 portsweep. 1
```

族群1只有一个数据点

K值的选择

k的值到底设置为多少比较合适呢？
在示例的样本数据中，有23个类别，那么意味着k至少等于23
通常情况下我们要通过多次尝试才能找到最好的k值

那么什么样的k值才是“最好”的呢？
如果每个数据点都紧靠最近的质心，那么这个聚类是较优的

为了能够判断每个数据点到质心的距离，我们可以编写以下函数来判断：



```
1  /**
2   * 计算两个向量之间的距离
3   *
4   * @param a 向量1
5   * @param b 向量2
6   *      欧式距离:空间上两个点的距离=两个向量相应元素的差的平方和的平方根
7   */
8  def distance(a: Vector, b: Vector) = {
9      //求平方根
10     math.sqrt(
11         //将两个向量合并
12         a.toArray.zip(b.toArray)
13         //两个向量中的每个值相减
14         .map(d => d._1 - d._2)
15         //相间的值平方
16         .map(d => d * d)
17         //之后相加
18         .sum)
19 }
20
21 /**
22 * 计算数据点到聚类中心质心的距离
23 *
24 * @param datum 数据点
25 * @param model kmeans模型
26 */
27 def distToCentroid(datum: Vector, model: KMeansModel) = {
28     //得到该数据点的聚类中心
29     val cluster = model.predict(datum)
30     //得到该聚类中心的质心
31     val centroid = model.clusterCenters(cluster)
32     //计算距离
33     distance(centroid, datum)
34 }
35
36 /**
37 * 根据各个数据点到该数据点聚类中心质心的距离来判断该模型优劣
38 * @param data 样本数据
39 * @param k k值
40 */
41 def clusteringScore(data: RDD[Vector], k: Int) = {
42     val kmeans = new KMeans()
43     //设置k值
44     kmeans.setK(k)
45     val model = kmeans.run(data)
46     //计算样本数据到其各自质心的记录的平均值
47     data.map { datum =>
48         distToCentroid(datum, model)
49     }.mean()
50 }
```

有了判断模型优劣的标准之后，就可以通过取不同的k值来观察模型：

```
1 //取不同k值观察模型优劣
2 (5 to 40 by 5).map { k =>
3     (k, clusteringScore(data, k))
4 }.foreach(println)
```

输出结果如下：



```
1 (5, 1779. 3473960726312)
2 (10, 1054. 5660956505587)
3 (15, 998. 6026754769782)
4 (20, 438. 0714456623944)
5 (25, 386. 70458251397577)
6 (30, 329. 4472646112194)
7 (35, 644. 939843705805)
8 (40, 221. 3547720824891)
```

可以看到，平均距离随着k的增大而降低

这点是毫无疑问的，因为随着族群点的增加，数据点离最近的质心肯定更近，当族群点等于数据量的时候平均距离为0，每个数据点都是自己构成的质心

而一个很奇怪的现象是，k=35时的距离竟然比k=30的时候要大

这是因为KMeans的迭代过程是从一个随机点开始的，因此可能收敛于一个局部最优解

k=35的情况可能是随机初始的质心造成的，也可能是由于算法在达到局部最小值之前就结束了

为了解决这个可能存在的问题，我们可以通过对固定的k值多次聚类，每次都随机不同的初始质心，然后在其中选择最优的

Spark Mllib提供了设置KMeans运行次数的方法，在clusteringScore函数中加入：

```
1 //设置该k值的聚类次数
2 kmeans.setRuns(10)
3 //设置迭代过程中, 质心的最小移动值, 默认为1.0e-4
4 kmeans.setEpsilon(1.0e-6)
```

setEpsilon设置迭代过程中,质心的最小移动值，移动值越小使得迭代的时间越多，聚类的结果更优化

现在我们重新选择k值进行评估：

```
1 (50 to 130 by 10).map { k =>
2   (k, clusteringScore(data, k))
3 }.foreach(println)
```

输出结果：

```
1 (50, 186. 72154668205906)
2 (60, 144. 68385906795461)
3 (70, 119. 93958735623515)
4 (80, 109. 46520683932486)
5 (90, 94. 08809036196241)
6 (100, 76. 37377878951273)
7 (110, 74. 78271035271912)
8 (120, 71. 515517236215)
9 (130, 65. 13241545688685)
```

这个时候随着k的增大，平均距离在持续减小

但是这个减少的幅度是有临界点的，当k值超过这个临界点，即使继续增大，也不会显著地降低距离

这个临界点就是我们要找的最好的k值

从输出结果中可以看到，k值应该取130

我们再次打印出每个族群包含的类别和个数信息：



```
1 0 neptune. 48517
2 0 nmap. 93
3 0 normal. 2974
4 0 portsweep. 904
5 0 rootkit. 3
6 0 satan. 222
7 0 teardrop. 865
8 0 warezmaster. 1
9 1 portsweep. 1
10 2 warezclient. 59
11 3 multihop. 1
12 3 normal. 1
13 4 normal. 1
14 5 normal. 310
15 6 normal. 1
16 6 warezmaster. 15
17 7 normal. 22
18 8 normal. 713
19 9 normal. 19
20 10 normal. 1
21 11 normal. 6
22 12 normal. 1
23 13 normal. 2
24 14 back. 2155
25 15 normal. 1
26 16 normal. 102
27 17 normal. 17
28 18 smurf. 227840
29 19 normal. 1
30 ...
31 ...
32 117 warezmaster. 1
33 118 normal. 3
34 119 normal. 2
35 120 back. 9
36 120 normal. 23
37 121 back. 1
38 121 normal. 362
39 122 multihop. 1
40 122 normal. 5174
41 122 smurf. 177
42 122 warezclient. 31
43 123 normal. 1242
44 124 normal. 351
45 125 back. 18
46 125 normal. 5
47 126 normal. 2
48 127 normal. 1
49 128 multihop. 1
50 128 normal. 376
51 128 rootkit. 1
52 129 normal. 7
```

结果比第一次好上很多了

现在可以使用这个模型对全体数据进行聚类了，使用这个模型可以将数据中离质心最远的点找出来

将这个点到质心的距离设置为阈值

当有新的数据进来时，判断这个数据到其质心的距离是否超过这个阈值

超过就发出警报进行异常检车

总结

在本篇中，样本数据的类别型特征被直接跳过，但是在实际场景中是不能这么做的
正确的做法应该是讲类别型特征转换为数值型特征来训练，得到的模型和之前讨论过的将不太一样
但是模型模型的训练和寻找最优k值的过程是一致的

KMeans训练出来的模型还可以和Spark Streaming相结合，搭建出实时的网络流量异常预警系统

[Github源码地址](#)

作者：[@小黑](#)

顶 0 踩 0

- [上一篇](#) MapReduce的类型与格式
- [下一篇](#) MapReduce高级特性

相关文章推荐

- 网络异常模拟测试方法
- MySQL在微信支付下的高可用运营--莫晓东
- 02storm聚类尝试kmeans
- 容器技术在58同城的实践--姚远
- 基本Kmeans算法介绍及其实现
- SDCC 2017之容器技术实战线上峰会
- Kmeans、Kmeans++和KNN算法比较
- SDCC 2017之数据库技术实战线上峰会
- KMeans的C++及Python实现
- 腾讯云服务器架构实现介绍--董晓杰
- KMeans聚类算法思想与可视化
- 微博热点事件背后的数据库运维心得--张冬洪
- KMeans
- KMeans和KMedoid 的Matlab实现
- 基于Kmeans的证件照背景色替换算法
- 机器学习算法与Python实践之（五）k均值聚类（...



迷你仓



人脸识别



口红排名



赴美生子费用



人脸识别算法

查看评论



qq_22245863

1楼 2017-02-08 10:37发表

val label = buffer.remove(buffer.length - 1) ? remove是您自己写的函数吗？

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场



可视化数据



人脸识别算法

