# 日志分析实战之清洗、统计网站信息小教程

作者：pig2
时间：20171017
出处：about 云

**目录**

# 1.使用 spark&Scala 分析 Apache 日志

**问题导读**

**1.如何进入 spark shell？**
**2.spark shell 中如何加载外部文件？**
**3.spark 中读取文件后做了哪些操作？**

about 云日志分析，那么过滤清洗日志。该如何实现。这里参考国外的一篇文章，总结分享给大家。
使用 spark 分析网站访问日志，日志文件包含数十亿行。现在开始研究 spark 使用，他是如何工作的。几年前使用 hadoop，后来发现 spark 也是容易的。
下面是需要注意的：

如果你已经知道如何使用 spark 并想知道如何处理 spark 访问日志记录，我写了这篇短的文章，介绍如何从 Apache 访问日志文件中生成 URL 点击率的排序

# 安装

spark 安装需要安装 hadoop，并且二者版本要合适。安装可参考下面文章

about 云日志分析项目准备 6：Hadoop、Spark 集群搭建

http://www.aboutyun.com/forum.php?mod=viewthread&tid=20620

进入

[Bash shell] 纯文本查看 复制代码

? 
1    ./bin/spark-shell

可能会出错

[Bash shell] 纯文本查看 复制代码

? 
1    java.io.FileNotFoundException: File
     file:/data/spark_data/history/event-log does not exist

解决办法：

[Bash shell] 纯文本查看 复制代码

? 
1    mkdir -p /data/spark_data/history/event-log


详细错误如下

[Bash shell] 纯文本查看 复制代码

? 
001 17/10/08 17:00:23 INFO client.AppClient$ClientEndpoint: Executor updated: app-201
002 17/10/08 17:00:25 ERROR spark.SparkContext: Error initializing SparkContext.
003 java.io.FileNotFoundException: File file:/data/spark_data/history/event-log does
004          at org.apache.hadoop.fs.RawLocalFileSystem.deprecatedGetFileStatu
005          at org.apache.hadoop.fs.RawLocalFileSystem.getFileLinkStatusInter
006          at org.apache.hadoop.fs.RawLocalFileSystem.getFileStatus(RawLocal
007          at org.apache.hadoop.fs.FilterFileSystem.getFileStatus(FilterFile
008          at org.apache.spark.scheduler.EventLoggingListener.start(EventLog
009          at org.apache.spark.SparkContext.<init>(SparkContext.scala:549)
010          at org.apache.spark.repl.SparkILoop.createSparkContext(SparkILoop
011          at $line3.$read$$iwC$$iwC.<init>(<console>:15)
012          at $line3.$read$$iwC.<init>(<console>:24)
013          at $line3.$read.<init>(<console>:26)
014          at $line3.$read$.<init>(<console>:30)
015          at $line3.$read$.<clinit>(<console>)
016          at $line3.$eval$.<init>(<console>:7)
017          at $line3.$eval$.<clinit>(<console>)

```
018                    at $line3.$eval.$print(<console>)
019                    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
020                    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
021                    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
022                    at java.lang.reflect.Method.invoke(Method.java:497)
023                    at org.apache.spark.repl.SparkIMain$ReadEvalPrint.call(SparkIMain
024                    at org.apache.spark.repl.SparkIMain$Request.loadAndRun(SparkIMain
025                    at org.apache.spark.repl.SparkIMain.loadAndRunReq$1(SparkIMain.sc
026                    at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:87
027                    at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:81
028                    at org.apache.spark.repl.SparkILoop.reallyInterpret$1(SparkILoop.
029                    at org.apache.spark.repl.SparkILoop.interpretStartingWith(SparkIL
030                    at org.apache.spark.repl.SparkILoop.command(SparkILoop.scala:814)
031                    at org.apache.spark.repl.SparkILoopInit$$anonfun$initializeSpark$
032                    at org.apache.spark.repl.SparkILoopInit$$anonfun$initializeSpark$
033                    at org.apache.spark.repl.SparkIMain.beQuietDuring(SparkIMain.scal
034                    at org.apache.spark.repl.SparkILoopInit$class.initializeSpark(Spa
035                    at org.apache.spark.repl.SparkILoop.initializeSpark(SparkILoop.sc
036                    at
037 org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$repl$SparkILoop$$proce
038                    at org.apache.spark.repl.SparkILoopInit$class.runThunks(SparkILoo
039                    at org.apache.spark.repl.SparkILoop.runThunks(SparkILoop.scala:64
040                    at org.apache.spark.repl.SparkILoopInit$class.postInitialization(
041                    at org.apache.spark.repl.SparkILoop.postInitialization(SparkILoop
042                    at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
043                    at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
044                    at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
045                    at scala.tools.nsc.util.ScalaClassLoader$.savingContextLoader(Sca
046                    at org.apache.spark.repl.SparkILoop.org$apache$spark$repl$SparkIL
047                    at org.apache.spark.repl.SparkILoop.process(SparkILoop.scala:1059
048                    at org.apache.spark.repl.Main$.main(Main.scala:31)
049                    at org.apache.spark.repl.Main.main(Main.scala)
050                    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
051                    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
052                    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
053                    at java.lang.reflect.Method.invoke(Method.java:497)
054                    at org.apache.spark.deploy.SparkSubmit$.org$apache$spark$deploy$S
055                    at org.apache.spark.deploy.SparkSubmit$.doRunMain$1(SparkSubmit.s
056                    at org.apache.spark.deploy.SparkSubmit$.submit(SparkSubmit.scala:
057                    at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:12
058                    at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
059
060                    at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
061
```

```
062 java.lang.NullPointerException
063                 at org.apache.spark.sql.SQLContext$.createListenerAndUI(SQLContex
064                 at org.apache.spark.sql.hive.HiveContext.<init>(HiveContext.scala
065                 at sun.reflect.NativeConstructorAccessorImpl.newInstance0(Native
066                 at sun.reflect.NativeConstructorAccessorImpl.newInstance(NativeCo
067                 at sun.reflect.DelegatingConstructorAccessorImpl.newInstance(Dele
068                 at java.lang.reflect.Constructor.newInstance(Constructor.java:422
069                 at org.apache.spark.repl.SparkILoop.createSQLContext(SparkILoop.s
070                 at $iwC$$iwC.<init>(<console>:15)
071                 at $iwC.<init>(<console>:24)
072                 at <init>(<console>:26)
073                 at .<init>(<console>:30)
074                 at .<clinit>(<console>)
075                 at .<init>(<console>:7)
076                 at .<clinit>(<console>)
077                 at $print(<console>)
078                 at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
079                 at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
080                 at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
081                 at java.lang.reflect.Method.invoke(Method.java:497)
082                 at org.apache.spark.repl.SparkIMain$ReadEvalPrint.call(SparkIMain
083                 at org.apache.spark.repl.SparkIMain$Request.loadAndRun(SparkIMain
084                 at org.apache.spark.repl.SparkIMain.loadAndRunReq$1(SparkIMain.sc
085                 at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:87
086                 at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:81
087                 at org.apache.spark.repl.SparkILoop.reallyInterpret$1(SparkILoop.
088                 at org.apache.spark.repl.SparkILoop.interpretStartingWith(SparkIL
089                 at org.apache.spark.repl.SparkILoop.command(SparkILoop.scala:814)
090                 at org.apache.spark.repl.SparkILoopInit$$anonfun$initializeSpark$
091                 at org.apache.spark.repl.SparkILoopInit$$anonfun$initializeSpark$
092                 at org.apache.spark.repl.SparkIMain.beQuietDuring(SparkIMain.scal
093                 at org.apache.spark.repl.SparkILoopInit$class.initializeSpark(Spa
094                 at org.apache.spark.repl.SparkILoop.initializeSpark(SparkILoop.sc
095                 at
096 org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$repl$SparkILoop$$proce
097                 at org.apache.spark.repl.SparkILoopInit$class.runThunks(SparkILoo
098                 at org.apache.spark.repl.SparkILoop.runThunks(SparkILoop.scala:64
099                 at org.apache.spark.repl.SparkILoopInit$class.postInitialization(
100                 at org.apache.spark.repl.SparkILoop.postInitialization(SparkILoop
101                 at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
102                 at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
103                 at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
104                 at scala.tools.nsc.util.ScalaClassLoader$.savingContextLoader(Sca
105                 at org.apache.spark.repl.SparkILoop.org$apache$spark$repl$SparkIL
```

```
106              at org.apache.spark.repl.SparkILoop.process(SparkILoop.scala:1059
107              at org.apache.spark.repl.Main$.main(Main.scala:31)
108              at org.apache.spark.repl.Main.main(Main.scala)
109              at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
110              at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
111              at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
112              at java.lang.reflect.Method.invoke(Method.java:497)
113              at org.apache.spark.deploy.SparkSubmit$.org$apache$spark$deploy$S
114              at org.apache.spark.deploy.SparkSubmit$.doRunMain$1(SparkSubmit.s
115              at org.apache.spark.deploy.SparkSubmit$.submit(SparkSubmit.scala:
116              at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:12
117              at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
118
119<console>:16: error: not found: value sqlContext
120              import sqlContext.implicits._
121                                 ^
122<console>:16: error: not found: value sqlContext
              import sqlContext.sql
                                 ^
```

进入 spark shell



[Bash shell] 纯文本查看 复制代码

**?**

**1**    val textFile=sc.textFile("file:///data/spark/README.md")

说明：

记得这里如果自己创建的文件可能会读取不到。报错如下

## [Bash shell] 纯文本查看 复制代码

**?**

```
001 java.io.FileNotFoundException: File file:/data/spark/change.txt does not exist
002              at org.apache.hadoop.fs.RawLocalFileSystem.deprecatedGetFileStatu
003              at org.apache.hadoop.fs.RawLocalFileSystem.getFileLinkStatusInter
004              at org.apache.hadoop.fs.RawLocalFileSystem.getFileStatus(RawLocal
005              at org.apache.hadoop.fs.FilterFileSystem.getFileStatus(FilterFile
006              at org.apache.hadoop.fs.ChecksumFileSystem$ChecksumFSInputChecker
007              at org.apache.hadoop.fs.ChecksumFileSystem.open(ChecksumFileSyste
008              at org.apache.hadoop.fs.FileSystem.open(FileSystem.java:766)
009              at org.apache.hadoop.mapred.LineRecordReader.<init>(LineRecordRea
010              at org.apache.hadoop.mapred.TextInputFormat.getRecordReader(TextI
011              at org.apache.spark.rdd.HadoopRDD$$anon$1.<init>(HadoopRDD.scala:
012              at org.apache.spark.rdd.HadoopRDD.compute(HadoopRDD.scala:211)
013              at org.apache.spark.rdd.HadoopRDD.compute(HadoopRDD.scala:101)
014              at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:306
015              at org.apache.spark.rdd.RDD.iterator(RDD.scala:270)
016              at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD
017              at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:306
018              at org.apache.spark.rdd.RDD.iterator(RDD.scala:270)
019              at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.scala
020              at org.apache.spark.scheduler.Task.run(Task.scala:89)
021              at org.apache.spark.executor.Executor$TaskRunner.run(Executor.sca
022              at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolEx
023              at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolE
024              at java.lang.Thread.run(Thread.java:745)
025
026 Driver stacktrace:
027              at
028 org.apache.spark.scheduler.DAGScheduler.org$apache$spark$scheduler$DAGScheduler$$
029              at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.
030              at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.
031              at scala.collection.mutable.ResizableArray$class.foreach(Resizabl
032              at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.scala
033              at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGSchedule
034              at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSet
035              at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSet
036              at scala.Option.foreach(Option.scala:236)
```

```
037              at org.apache.spark.scheduler.DAGScheduler.handleTaskSetFailed(DA
038              at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.doOnRe
039              at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRece
040              at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRece
041              at org.apache.spark.util.EventLoop$$anon$1.run(EventLoop.scala:48
042              at org.apache.spark.scheduler.DAGScheduler.runJob(DAGScheduler.sc
043              at org.apache.spark.SparkContext.runJob(SparkContext.scala:1832)
044              at org.apache.spark.SparkContext.runJob(SparkContext.scala:1845)
045              at org.apache.spark.SparkContext.runJob(SparkContext.scala:1858)
046              at org.apache.spark.SparkContext.runJob(SparkContext.scala:1929)
047              at org.apache.spark.rdd.RDD.count(RDD.scala:1157)
048              at $iwC$$iwC$$iwC$$iwC$$iwC$$iwC$$iwC$$iwC.<init>(<console>:30)
049              at $iwC$$iwC$$iwC$$iwC$$iwC$$iwC$$iwC.<init>(<console>:35)
050              at $iwC$$iwC$$iwC$$iwC$$iwC$$iwC.<init>(<console>:37)
051              at $iwC$$iwC$$iwC$$iwC$$iwC.<init>(<console>:39)
052              at $iwC$$iwC$$iwC$$iwC.<init>(<console>:41)
053              at $iwC$$iwC$$iwC.<init>(<console>:43)
054              at $iwC$$iwC.<init>(<console>:45)
055              at $iwC.<init>(<console>:47)
056              at <init>(<console>:49)
057              at .<init>(<console>:53)
058              at .<clinit>(<console>)
059              at .<init>(<console>:7)
060              at .<clinit>(<console>)
061              at $print(<console>)
062              at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
063              at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
064              at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
065              at java.lang.reflect.Method.invoke(Method.java:497)
066              at org.apache.spark.repl.SparkIMain$ReadEvalPrint.call(SparkIMain
067              at org.apache.spark.repl.SparkIMain$Request.loadAndRun(SparkIMain
068              at org.apache.spark.repl.SparkIMain.loadAndRunReq$1(SparkIMain.sc
069              at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:87
070              at org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:81
071              at org.apache.spark.repl.SparkILoop.reallyInterpret$1(SparkILoop.
072              at org.apache.spark.repl.SparkILoop.interpretStartingWith(SparkIL
073              at org.apache.spark.repl.SparkILoop.command(SparkILoop.scala:814)
074              at org.apache.spark.repl.SparkILoop.processLine$1(SparkILoop.scal
075              at org.apache.spark.repl.SparkILoop.innerLoop$1(SparkILoop.scala:
076              at org.apache.spark.repl.SparkILoop.org$apache$spark$repl$SparkIL
077              at
078 org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$repl$SparkILoop$$proce
079              at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
080              at org.apache.spark.repl.SparkILoop$$anonfun$org$apache$spark$rep
```

```
081            at scala.tools.nsc.util.ScalaClassLoader$.savingContextLoader(Sca
082            at org.apache.spark.repl.SparkILoop.org$apache$spark$repl$SparkII
083            at org.apache.spark.repl.SparkILoop.process(SparkILoop.scala:105
084            at org.apache.spark.repl.Main$.main(Main.scala:31)
085            at org.apache.spark.repl.Main.main(Main.scala)
086            at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
087            at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccess
088            at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMeth
089            at java.lang.reflect.Method.invoke(Method.java:497)
090            at org.apache.spark.deploy.SparkSubmit$.org$apache$spark$deploy$S
091            at org.apache.spark.deploy.SparkSubmit$.doRunMain$1(SparkSubmit.s
092            at org.apache.spark.deploy.SparkSubmit$.submit(SparkSubmit.scala:
093            at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:12
094            at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
095Caused by: java.io.FileNotFoundException: File file:/data/spark/change.txt does r
096            at org.apache.hadoop.fs.RawLocalFileSystem.deprecatedGetFileStatu
097            at org.apache.hadoop.fs.RawLocalFileSystem.getFileLinkStatusInter
098            at org.apache.hadoop.fs.RawLocalFileSystem.getFileStatus(RawLocal
099            at org.apache.hadoop.fs.FilterFileSystem.getFileStatus(FilterFile
100            at org.apache.hadoop.fs.ChecksumFileSystem$ChecksumFSInputChecker
101            at org.apache.hadoop.fs.ChecksumFileSystem.open(ChecksumFileSyste
102            at org.apache.hadoop.fs.FileSystem.open(FileSystem.java:766)
103            at org.apache.hadoop.mapred.LineRecordReader.<init>(LineRecordRea
104            at org.apache.hadoop.mapred.TextInputFormat.getRecordReader(TextI
105            at org.apache.spark.rdd.HadoopRDD$$anon$1.<init>(HadoopRDD.scala:
106            at org.apache.spark.rdd.HadoopRDD.compute(HadoopRDD.scala:211)
107            at org.apache.spark.rdd.HadoopRDD.compute(HadoopRDD.scala:101)
108            at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:306
109            at org.apache.spark.rdd.RDD.iterator(RDD.scala:270)
110            at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDI
111            at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:306
112            at org.apache.spark.rdd.RDD.iterator(RDD.scala:270)
113            at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.scala
114            at org.apache.spark.scheduler.Task.run(Task.scala:89)
115            at org.apache.spark.executor.Executor$TaskRunner.run(Executor.sca
116            at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolEx
               at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolE
               at java.lang.Thread.run(Thread.java:745)
```

需要是文件权限为 500,才可以读取到。

执行

## [Bash shell] 纯文本查看 复制代码

?
1    textFile.count



## [Bash shell] 纯文本查看 复制代码

?

1    textFile.first

输出如下内容

## [Bash shell] 纯文本查看 复制代码

?

```
scala> textFile.first
17/10/08 18:34:23 INFO spark.SparkContext: Starting job: first at
<console>:30
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Got job 1 (first at
<console>:30) with 1 output partitions
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Final stage: ResultStage 1
(first at <console>:30)
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Parents of final stage:
List()
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Missing parents: List()
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Submitting ResultStage 1
(file:///data/spark/README.md MapPartitionsRDD[1] at textFile at
<console>:27), which has no missing parents
17/10/08 18:34:23 INFO storage.MemoryStore: Block broadcast_2 stored as
values in memory (estimated size 3.1 KB, free 517.2 MB)
17/10/08 18:34:23 INFO storage.MemoryStore: Block broadcast_2_piece0
stored as bytes in memory (estimated size 1843.0 B, free 517.2 MB)
17/10/08 18:34:23 INFO storage.BlockManagerInfo: Added broadcast_2_piece0
in memory on 192.168.1.10:41717 (size: 1843.0 B, free: 517.4 MB)
17/10/08 18:34:23 INFO spark.SparkContext: Created broadcast 2 from
broadcast at DAGScheduler.scala:1006
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Submitting 1 missing tasks
from ResultStage 1 (file:///data/spark/README.md MapPartitionsRDD[1] at
textFile at <console>:27)
17/10/08 18:34:23 INFO scheduler.TaskSchedulerImpl: Adding task set 1.0
with 1 tasks
17/10/08 18:34:23 INFO scheduler.TaskSetManager: Starting task 0.0 in stage
1.0 (TID 2, slave2, partition 0,PROCESS_LOCAL, 2128 bytes)
17/10/08 18:34:23 INFO storage.BlockManagerInfo: Added broadcast_2_piece0
in memory on slave2:35228 (size: 1843.0 B, free: 517.4 MB)
17/10/08 18:34:23 INFO scheduler.TaskSetManager: Finished task 0.0 in stage
1.0 (TID 2) in 116 ms on slave2 (1/1)
17/10/08 18:34:23 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 1.0,
whose tasks have all completed, from pool
17/10/08 18:34:23 INFO scheduler.DAGScheduler: ResultStage 1 (first at
<console>:30) finished in 0.117 s
17/10/08 18:34:23 INFO scheduler.DAGScheduler: Job 1 finished: first at
```

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20

```
<console>:30, took 0.161753 s
res1: String = # Apache Spark
```

# 2.导入日志清洗代码并打包

**问题导读**

**1.通过什么菜单项可以导入源码？**
**2.打 jar 包需要哪些步骤？**
**3.如何找到 jar 生成路径？**

上一篇：
about 云日志分析实战之清洗日志 1：使用 spark&Scala 分析 Apache 日志
http://www.aboutyun.com/forum.php?mod=viewthread&tid=22856

前面测试了一下 spark，准备好环境，下面开始动工源码。分析清洗日志，这里面的代码还是比较复杂的。
对于 iis 日志，可参考
about 云日志分析项目准备 10-3：Spark Local 模式之 Log 文本清洗
http://www.aboutyun.com/forum.php?mod=viewthread&tid=21135
对于 Apache 日志，国外已经实现。源码 git 地址
https://github.com/alvinj/ScalaApacheAccessLogParser

网盘下载地址
链接：http://pan.baidu.com/s/1jIj87wM 密码：p0zd

这里从上面下载下来，然后导入 IntelliJ IDEA ，然后打包。

# 导入源码

首先 file-》open

选择源码文件

导入之后看到下面内容



对于 spark 环境不熟悉或则不会操作可参考
spark 开发环境详细教程 1：IntelliJ IDEA 使用详细说明
http://www.aboutyun.com/forum.php?mod=viewthread&tid=22320

spark 开发环境详细教程 3：IntelliJ IDEA 创建项目
http://www.aboutyun.com/forum.php?mod=viewthread&tid=22410

打包

上面我们准备了源码，然后将源码打成 jar 包，供我们项目中使用。
首先打开 project structure,

选择依赖

填写主类



点击确定

选择菜单 Build Artifacts

点击 build



最后生成 jar 包，在 terminal 中会显示输出 jar 包路径

找到生成 jar 包。我们就可以使用了。

# 3.如何在 spark shell 中导入自定义包

问题导读

**1.自定义包，本文放到哪个路径下面？**
**2.复制包之后，需要做哪些权限操作？**
**3.如何验证导入是否成功？**

上一篇
about 云日志分析实战之清洗日志 2：导入日志清洗代码并打包
http://www.aboutyun.com/forum.php?mod=viewthread&tid=22862

上一篇文章，生成了包，那么这个包该如何加载到 spark 环境中，并且为我们所使用。那么首先改如何加载这个包。
首先将这个包放到 spark 中的 lib 文件夹下。



在复制到 Linux 中，首先需要修改的就是权限。
我们看到用户和组的权限为 500，并且用户，所属组，及其它用户都为满权限，
可以通过下面命令来实现

## [Bash shell] 纯文本查看 复制代码

?
**1**    sudo chown 500:500 ScalaApacheAccessLogParser-master.jar

## [Bash shell] 纯文本查看 复制代码

?

1　　sudo chmod -R a+r ScalaApacheAccessLogParser-master.jar

## [Bash shell] 纯文本查看 复制代码

?

1　　sudo chmod -R a+w ScalaApacheAccessLogParser-master.jar

## [Bash shell] 纯文本查看 复制代码

?

1　　sudo chmod -R a+x ScalaApacheAccessLogParser-master.jar

通过上面命令即可实现授权。

授权完毕，接着我们就需要把这个包，加载到 spark shell 环境中。

## [Bash shell] 纯文本查看 复制代码

?

1　　./bin/spark-shell --jars lib/ScalaApacheAccessLogParser-master.jar

接着我们执行导入 jar 包

## [Bash shell] 纯文本查看 复制代码

?

1　　import com.alvinalexander.accesslogparser._



至此我们就可以使用第三方包了。

问题：

同时尝试了比较多的导入方式，没有成功，记录下来共大家借鉴。

## [Bash shell] 纯文本查看 复制代码

?

1

```
./bin/spark-shell - master spark://master:7077  - jars
ScalaApacheAccessLogParser-master.jar
```

## [Bash shell] 纯文本查看 复制代码

?

1

```
MASTER=local[4]
ADD_JARS=/data/spark/lib/AlsApacheLogParser.jar ./bin/spark-shell
```

导入的时候，并不会报错，但是 import 的时候，报错。



# 4.统计网站相关信息

**问题导读**

**1.如何统计网站总的点击量？**

**2.如何实现统计不能访问网页的个数？**

**3.文章中如何定义和使用 Scala 函数的？**

上一篇

about 云日志分析实战之清洗日志 3：如何在 spark shell 中导入自定义包

http://www.aboutyun.com/forum.php?mod=viewthread&tid=22881

上一篇，我们已经添加了清洗日志的核心代码，那么剩下的我们就可以统计相关信息，比如最简单的找到不能访问的网页。

导入之后，我们创建 AccessLogParser 实例

## [Bash shell] 纯文本查看 复制代码

?

**1**    `val p = new AccessLogParser`

这个很重要，在后面我们会用到
首先我们需要加载一部分日志样例。

## [Bash shell] 纯文本查看 复制代码

?

```
      192.168.169.50 - - [17/Feb/2012:10:09:13 +0800] "GET /favicon.ico HTTP/1.1"
      404 288 "-" "360se"
      192.168.169.50 - - [17/Feb/2012:10:36:26 +0800] "GET / HTTP/1.1" 403 5043
      "-" "Mozilla/5.0 (Windows NT 5.1; rv:6.0) Gecko/20100101 Firefox/6.0"
      192.168.169.50 - - [17/Feb/2012:10:36:26 +0800] "GET
      /icons/powered_by_rh.png HTTP/1.1" 200 1213 "http://192.168.55.230/"
      "Mozilla/5.0 (Windows NT 5.1; rv:6.0) Gecko/20100101 Firefox/6.0"
      192.168.169.50 - - [17/Feb/2012:10:09:10 +0800] "GET
      /icons/powered_by_rh.png HTTP/1.1" 200 1213 "http://192.168.55.230/"
01    "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1; Trident/4.0;
02    InfoPath.2; 360SE)"
03    192.168.55.230 - - [24/Feb/2012:09:48:58 +0800] "GET /favicon.ico HTTP/1.1"
04    404 288 "-" "Mozilla/5.0 (X11; U; Linux x86_64; en-US; rv:1.9.2.24)
05    Gecko/20111109 CentOS/3.6-3.el5.centos Firefox/3.6.24"
06    192.168.169.50 - - [24/Feb/2012:09:45:03 +0800] "GET /server-status
07    HTTP/1.1" 404 290 "-" "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1;
08    Trident/4.0; InfoPath.2; 360SE)"
09    192.168.55.230 - - [24/Feb/2012:09:49:02 +0800] "GET / HTTP/1.1" 403 5043
10    "-" "Mozilla/5.0 (X11; U; Linux x86_64; en-US; rv:1.9.2.24) Gecko/20111109
      CentOS/3.6-3.el5.centos Firefox/3.6.24"
      192.168.55.230 - - [24/Feb/2012:09:49:02 +0800] "GET /icons/apache_pb.gif
      HTTP/1.1" 200 2326 "http://192.168.55.230/" "Mozilla/5.0 (X11; U; Linux
      x86_64; en-US; rv:1.9.2.24) Gecko/20111109 CentOS/3.6-3.el5.centos
      Firefox/3.6.24"
      192.168.55.230 - - [24/Feb/2012:09:49:02 +0800] "GET
      /icons/powered_by_rh.png HTTP/1.1" 200 1213 "http://192.168.55.230/"
      "Mozilla/5.0 (X11; U; Linux x86_64; en-US; rv:1.9.2.24) Gecko/20111109
      CentOS/3.6-3.el5.centos Firefox/3.6.24"
```

```
192.168.55.230 - - [24/Feb/2012:09:49:20 +0800] "GET /server-status
HTTP/1.1" 404 290 "-" "Mozilla/5.0 (X11; U; Linux x86_64; en-US;
rv:1.9.2.24) Gecko/20111109 CentOS/3.6-3.el5.centos Firefox/3.6.24"
```

将其保存为 aboutyun.log

将其上传到 hadoop

## [Bash shell] 纯文本查看 复制代码

?
1
```
hadoop fs -put aboutyun.log /
```

上传成功验证



# 统计网站总的点击量

接着我们加载文件。

## [Bash shell] 纯文本查看 复制代码

?
1
```
var log=sc.textFile("/aboutyun.log")
```

这里 sc 是系统已经初始化的，我们可以直接使用，可以理解为 sparkContext 的实例



加载之后，我们统计行数，也可以理解为统计网站总的点击量。这时候我们就看到总点击量为 10

# 统计网站不能访问网页的数量

首先我们定义一个函数，获取一条记录的 httpStatusCode，也就是返回码

[Scala] 纯文本查看 复制代码

```
?
1    def getStatusCode(line: Option[AccessLogRecord]) = {
2        line match {
3            case Some(l) => l.httpStatusCode
4            case None => "0"
5        }
6    }
```



定义函数之后，我们接着使用

[Bash shell] 纯文本查看 复制代码

?

1     log.filter(line => getStatusCode(p.parseRecord(line)) == "404").count

上面的 p 是我们前面定义的对象。

val p = new AccessLogParser，然后调用了 parseRecord 函数。这些都是 jar 包的内容。大家可以详细看看。



这样 404 网页的个数就统计出来了。后面我们可以做一些更加复杂的内容

#################

**补充说明**

**1.**在统计日志测试的时候，文件一定标准，否则会统计错误，比如日志要换行

**2.函数定义**

附上所用函数的相关信息

**Option and Either**

Option 和 Either 都是用来让返回值可以有两个选择

而 Option 是比较简单的版本，两个选择，一定是成功 Some，和失败 None

Option 意味着可能有值 some(x)，也可能没有值(用 None 对象，表示缺

失)，典型的例子就是从字典里取值

[Scala] 纯文本查看 复制代码

?

1     val capitals = Map("France" -> "Paris", "Japan" -> "Tokyo")

2     def show(x: Option[String]) = x match { //Option 类型，可选的 String

```
3            case Some(s) => s
4            case None => "?"
5    }
6    scala> show(capitals get "France")
7    res24: String = Paris
8    scala> show(capitals get "North Pole")
9    res25: String = ?
```

以前的方式，比如 Java，通过 null 来表示没有取到值，但是有的时候

null 可能作为合法值出现，就需要特殊处理，很麻烦

而 Scala 提供 option 来比较优雅的解决这个问题

Either，更为通用一些，可用自己定义两种选择，直接看个spark 源码中

的例子，

对于 PutResult 中的 data，有可能是 ByteBuffer 或者 Iterator

而使用的时候，使用 Left 和 Right 来选择到底用哪一个

[Scala] 纯文本查看 复制代码

```
?
     private[spark] case class PutResult(size: Long, data: Either[Iterator[_],
1    ByteBuffer])
2
3    PutResult(sizeEstimate, Left(values.iterator))PutResult(bytes.limit(),
     Right(bytes.duplicate()))
```

这里无论 option 或 either 都提高了极好的灵活性，在 Java 中如果要返

回一个有两种可能性的值就比较不那么优雅了

来自：

http://www.cnblogs.com/fxjwind/p/3338829.html

# 5.实现获取不能访问 url

**问题导读**

**1.在 url 中，如何过滤不需要的内容？**
**2.如何获取 404 记录并且获取字段？**
**3.获取不能访问 url 列表的思路是什么？**

about 云日志分析实战之清洗日志 4：统计网站相关信息

http://www.aboutyun.com/forum.php?mod=viewthread&tid=22900

上篇文章简单的统计了一些信息，下一步希望找到 404 对应的 url。

**思路：**

1.获取 request 字段

2.过滤不需要字符

3.实现获取 url，并打印输出

**1.创建 getRequest 函数获取 request 字段**

## [Scala] 纯文本查看 复制代码

```
?
1    // get the `request` field from an access log record
2    def getRequest(rawAccessLogString: String): Option[String] = {
3        val accessLogRecordOption = p.parseRecord(rawAccessLogString)
4        accessLogRecordOption match {
5            case Some(rec) => Some(rec.request)
6            case None => None
7        }
8    }
```

**2.创建 extractUriFromRequest 函数**

## [Scala] 纯文本查看 复制代码

```
?
1    // val request = "GET /foo HTTP/1.0"
```

2    def extractUriFromRequest(requestField: String) = requestField.split("
     ")(1)

这个目的大家可以猜猜它的作用

获取 404 页面，并且打印出请求的 URL.

## [Scala] 纯文本查看 复制代码

?

```
val distinctRecs = log.filter(line => getStatusCode(p.parseRecord(line))
== "404")
                                          .map(getRequest(_))
                                          .collect { case
Some(requestField) => requestField }
                                          .map(extractUriFromRequest(_))
                                          .distinct
```



## [Scala] 纯文本查看 复制代码

?

1    distinctRecs.count

## [Scala] 纯文本查看 复制代码

?

1    `distinctRecs.collect().foreach(println(_))`



**3.获取 url**

## [Scala] 纯文本查看 复制代码

?

```
1   val distinctRecs = log.filter(line => getStatusCode(p.parseRecord(line)) == "404")
2                                           .map(getRequest(_))
3                                           .collect { case
4   Some(requestField) => requestField }
5                                           .map(extractUriFromRequest(_))
                                            .distinct
```

通过上面看，其实挺简单。**Scala** 本身是非常简洁的。

**相关说明：**
上面看似简单，其实有很多需要说明的

val recs = log.filter(line => getStatusCode(p.parseRecord(line)) == "404").map(getRequest(_))

上面得出 404 对应的 url.getRequest 是上面我们定义的函数

val distinctRecs = log.filter(line => getStatusCode(p.parseRecord(line)) ==

"404").map(getRequest(_)).distinct

这里多了 distinct 是为了去重，下面是直接打印。

distinctRecs.collect().foreach(println(_))。

对于 extractUriFromRequest，这个主要为过滤我们不想要的内容。如下面，GET 和 HTTP/1.1 都不是我们想要的。所以我们取第二个元素即可。

## [Bash shell] 纯文本查看 复制代码

```
?
1    GET /foo HTTP/1.0
2    GET /foo HTTP/1.1
```

**知识补充：**

对于 collect（） 函数，是比较常见的，但是对于下面内容，是什么意思。

collect { case Some(requestField) => requestField }这个作用，类似 map。

##################

更多信息：

在 Scala 中，当我需要对集合的元素进行转换时，自然而然会使用到 map 方法。而当我们在对 tuple 类型的集合或者针对 Map 进行 map 操作时，通常更倾向于在 map 方法中使用 case 语句，这比直接使用_1 与_2 更加可读。例如：

## [Scala] 纯文本查看 复制代码

```
?
1    val languageToCount = Map("Scala" -> 10, "Java" -> 20, "Ruby" -> 5)
2    languageToCount map { case (_, count) => count + 1 }
```

然而对于上述场景，其实我们也可以使用 collect 方法：

## [Scala] 纯文本查看 复制代码

```
?
1    languageToCount collect { case (_, count) => count + 1 }
```

参考

http://www.jianshu.com/p/fa2ed7ed391e

# 6.获取 uri 点击量排序并得到最高的 url

**问题导读**

**1.读取日志的过程中，发生异常本文是如何解决的?**
**2.读取后，如何过滤异常的记录?**
**3.如何实现统计点击最高的记录?**

日志分析实战之清洗日志小实例 5：实现获取不能访问 url

http://www.aboutyun.com/forum.php?mod=viewthread&tid=22911

下面我们开始统计链接的点击量，并做排序。
我们统计记录的时候，为了防止空记录等异常的情况，我们创建一条空记录

## [Bash shell] 纯文本查看 复制代码

?

```
1    val nullObject = AccessLogRecord("", "", "", "", "GET /foo HTTP/1.1", "",
     "", "", "")
```

下面我们开始找点击量最高的链接。

**首先获取我们想要的 uri**

## [Scala] 纯文本查看 复制代码

?

```
1    val uriCounts = log.map(p.parseRecord(_).getOrElse(nullObject).request)
2                                    .map(_.split(" ")(1))
3                                    .filter(_ != "/foo")
```

上面的代码做一个简单解释:

p.parseRecord(_)解析记录

p.parseRecord(_).getOrElse(nullObject)如何没有取到值，则使用 nullObject，也就是我们上面定义的对象

p.parseRecord(_).getOrElse(nullObject).request 也就是我们取到 uri

.map(_.split(" ")(1))是取到我们过滤的 url，过滤掉不想要的版本等信息

  .filter(_ != "/foo")则是再次过滤掉/foo[也就是空记录]

这样就获取了 uri,然后我们输出

[Scala] 纯文本查看 复制代码

?

```
1    uriCounts.collect.foreach(print)
```



下面我们统计点击量

[Scala] 纯文本查看 复制代码

?

```
1    val uriCounts = log.map(p.parseRecord(_).getOrElse(nullObject).request)
2                                    .map(_.split(" ")(1))
3                                    .map(uri => (uri, 1))
4                                    .reduceByKey((a, b) => a + b)
```

**rdd** 转换为数组

[Scala] 纯文本查看 复制代码

?

```
1    val uriToCount = uriCounts.collect
```

**数组转换为序列并排序**

[Scala] 纯文本查看 复制代码

?

```
1    import scala.collection.immutable.ListMap
2    val uriHitCount = ListMap(uriToCount.toSeq.sortWith(_._2 > _._2):_*)
```



```
scala> val hitCount=ListMap(uriToCount.toSeq.sortWith(_._2>_._2):_*)
hitCount: scala.collection.immutable.ListMap[String,Int] = Map(/icons/powered_by_rh.png -> 3, /favicon.ico -> 2, / -> 2, /server-status -> 2, /icons/a
pache_pb.gif -> 1)

scala>
```

###############################

这里留下一个问题，如果上面元素不是 2，而是为 sortWith(_._1 > _._1)是对什么排序

## [Scala] 纯文本查看 复制代码

?
```
1    import scala.collection.immutable.ListMap
2    val uriHitCount = ListMap(uriToCount.toSeq.sortWith(_._1 > _._1):_*)
```
###############################



```
uriCount.take(1).foreach(println)

scala> val uriHitCount = ListMap(uriToCount.toSeq.sortWith(_._1 > _._1):_*)
uriHitCount: scala.collection.immutable.ListMap[String,Int] = Map(/server-status -> 2, /icons/powered_by_rh.png -> 3, /icons/apache_pb.gif -> 1, /favi
con.ico -> 2, / -> 2)

scala> val uriHitCount = ListMap(uriToCount.toSeq.sortWith(_._2 > _._2):_*)
uriHitCount: scala.collection.immutable.ListMap[String,Int] = Map(/icons/powered_by_rh.png -> 3, /favicon.ico -> 2, / -> 2, /server-status -> 2, /icon
s/apache_pb.gif -> 1)

scala>
```

输出

## [Scala] 纯文本查看 复制代码

?
```
1    uriHitCount.take(10).foreach(println)
```



```
scala> uriHitCount.take(10).foreach(println)
(/icons/powered_by_rh.png,3)
(/favicon.ico,2)
(/,2)
(/server-status,2)
(/icons/apache_pb.gif,1)

scala>
```

上面便是排序的结果

**点击最高的 uri**

如果想得出点击最高的 uri

## [Scala] 纯文本查看 复制代码

?

**1**     uriHitCount.take(1).foreach(println)



~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

知识补充：

Scala 代码看上去很少，但是内容却是很丰富的。上面用到的相关知识，这里补充，供大家能看懂上面代码

### getOrElse:

println(a.get("k1").getOrElse("default")) //根据 key 读取元素，不存在就替换成默认值

在 Spark 中写法是：persons.getOrElse("Spark",1000) //如果 persons 这个 Map 中包含有 Spark，取出它的值，如果没有，值就是 1000。

### reduce、reduceByKey

reduce(binary_function)

reduce 将 RDD 中元素前两个传给输入函数，产生一个新的 return 值，新产生的 return 值与 RDD 中下一个元素（第三个元素）组成两个元素，再被传给输入函数，直到最后只有一个值为止。

## [Scala] 纯文本查看 复制代码

?

**1**     val c = sc.parallelize(1 to 10)
**2**     c.reduce((x, y) => x + y)//结果 55

具体过程，RDD 有 1 2 3 4 5 6 7 8 9 10 个元素，

1+2=3
3+3=6
6+4=10
10+5=15
15+6=21
21+7=28
28+8=36
36+9=45

45+10=55

reduceByKey(binary_function)
reduceByKey 就是对元素为 KV 对的 RDD 中 Key 相同的元素的 Value 进行 binary_function 的 reduce 操作，因此，Key 相同

的多个元素的值被 reduce 为一个值，然后与原 RDD 中的 Key 组成一个新的 KV 对。
val a = sc.parallelize(List((1,2),(1,3),(3,4),(3,6)))
a.reduceByKey((x,y) => x + y).collect
//结果 Array((1,5), (3,10))


**Seq**

Sequence 都有一个预定义的顺序。
scala> Seq(1, 1, 2)
res3: Seq[Int] = List(1, 1, 2)
(注意返回的结果是一个 List。Seq 是一个 trait；List 是它的一个实现类。Seq 对象是一个工厂对象，正如你所看到

的，它会创建一个 List。)

**集合之间可以相互进行转换。**
def toArray : Array[A]
def toArray [B >: A] (implicit arg0: ClassManifest[B]) : Array[B]
def toBuffer [B >: A] : Buffer[B]
def toIndexedSeq [B >: A] : IndexedSeq[B]
def toIterable : Iterable[A]
def toIterator : Iterator[A]
def toList : List[A]
def toMap [T, U] (implicit ev: <:<[A, (T, U)]) : Map[T, U]
def toSeq : Seq[A]
def toSet [B >: A] : Set[B]
def toStream : Stream[A]
def toString () : String
def toTraversable : Traversable[A]

我们可以把一个 Map 转换成一个数组，然后得到一个键值对数组。

scala> Map(1 -> 2).toArray
res41: Array[(Int, Int)] = Array((1,2))

**sortWith**
排序操作（sorted, sortWith, sortBy）根据不同的条件对序列元素进行排序。

更多大家可以搜索


后面 about 云会有相关的日志实战视频，会通过 spark sql 等方式来实现。
相关文章链接：
http://www.aboutyun.com/forum.php?mod=group&fid=139
大家感兴趣可加入
日志分析实战群：238250736


更多：

搜索：微信号

wwwaboutyuncom



qq7 群：552029443 spark 学习
qq9 群：634068865 机器学习
qq5 群：432264021 storm 学习