

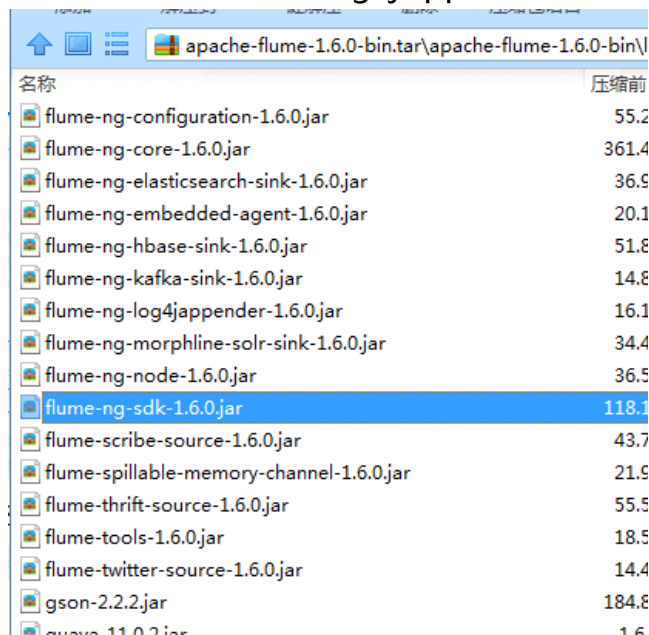
Flume_离线处理_日志收集

2019年3月10日 16:01

1. 发送日志到Flume

在日志服务器中，通过Log4jAppender将日志发往flume客户端

a. 在日志服务器应用中导入Log4jAppender相关开发包



b. 配置log4j配置文件，实现发送日志给flume

```
1 log4j.rootLogger = info,stdout,flume
2
3
4 log4j.appender.stdout = org.apache.log4j.ConsoleAppender
5 log4j.appender.stdout.Target = System.out
6 log4j.appender.stdout.layout = org.apache.log4j.PatternLayout
7 log4j.appender.stdout.layout.ConversionPattern = %m%n
8
9 log4j.appender.flume = org.apache.flume.clients.log4jappender.Log4jAppender
10 log4j.appender.flume.Hostname = hadoop01
11 log4j.appender.flume.Port = 44444
12 log4j.appender.stdout.layout = org.apache.log4j.PatternLayout
    log4j.appender.stdout.layout.ConversionPattern = %m%n
```

c. 在日志服务器的LogServlet中，通过log4j来发送日志

```
1 logger.info(line);
```

2. 开发客户端Agent

Hadoop03

```
1 #声明Agent
2 a1.sources = r1
3 a1.sinks = k1 k2
4 a1.channels = c1
5
6 #声明source
7 a1.sources.r1.type = avro
```

```

8   a1.sources.r1.bind = 0.0.0.0
9   a1.sources.r1.port = 44444
10
11  a1.sources.r1.interceptors = i1
12  a1.sources.r1.interceptors.i1.type = regex_extractor
13  a1.sources.r1.interceptors.i1.regex = ^(?:[^\]]*\]){14}\\d+\\d+_(\\d+)\\$.*$
14  a1.sources.r1.interceptors.i1.serializers = s1
15  a1.sources.r1.interceptors.i1.serializers.s1.name = timestamp
16
17  #声明sink
18  a1.sinks.k1.type = avro
19  a1.sinks.k1.hostname = hadoop01
20  a1.sinks.k1.port = 44444
21
22  a1.sinks.k2.type = avro
23  a1.sinks.k2.hostname = hadoop02
24  a1.sinks.k2.port = 44444
25
26  a1.sinkgroups = g1
27  a1.sinkgroups.g1.sinks = k1 k2
28  a1.sinkgroups.g1.processor.type = load_balance
29  a1.sinkgroups.g1.processor.backoff = true
30  a1.sinkgroups.g1.processor.selector = random
31
32  #声明channel
33  a1.channels.c1.type = memory
34  a1.channels.c1.capacity = 1000
35  a1.channels.c1.transactionCapacity = 100
36
37  #绑定关系
38  a1.sources.r1.channels = c1
39  a1.sinks.k1.channel = c1
40  a1.sinks.k2.channel = c1

```

3. 开发中心服务器Agent

hadoop01 hadoop02

```

1   #配置Agent
2   a1.sources = r1
3   a1.sinks = k1
4   a1.channels = c1
5
6   #声明Source
7   a1.sources.r1.type = avro
8   a1.sources.r1.bind = 0.0.0.0
9   a1.sources.r1.port = 44444
10
11  #声明sink
12  a1.sinks.k1.type = hdfs
13  a1.sinks.k1.hdfs.path = hdfs://hadoop01:9000/flux/reportTime=%Y-%m-%d
14  a1.sinks.k1.hdfs.rollInterval = 30
15  a1.sinks.k1.hdfs.rollSize = 0
16  a1.sinks.k1.hdfs.rollCount = 0
17  a1.sinks.k1.hdfs.fileType = DataStream
18  a1.sinks.k1.hdfs.timeZone = GMT+8
19
20  #声明channel

```

```

21  #声明channel
22  a1.channels.c1.type = memory
23  a1.channels.c1.capacity = 1000
24  a1.channels.c1.transactionCapacity = 100
25
26  #绑定关系
27  a1.sources.r1.channels = c1
    a1.sinks.k1.channel = c1

```

4. 遇到的问题

a. 找不到hadoop jar包

flume中的hdfs sink需要hadoop相关jar包的支持，

要么手动将hadoop相关jar包放置到flume的lib目录下

要么在本机中解压hadoop并将hadoop路径配置为HADOOP_HOME环境变量，使flume可以自动找到这些jar。

b. 产生大量小文件

hdfs sink的滚动条件设置不合理。

修改即可

```

1  a1.sinks.k1.hdfs.rollInterval = 30
2  a1.sinks.k1.hdfs.rollSize = 0
3  a1.sinks.k1.hdfs.rollCount = 0

```

c. 文件内容为乱码(序列化文件无法直接查看)

hdfs sink默认产生SequenceFile文件，无法直接查看

修改即可：

```

1  a1.sinks.k1.hdfs.fileType = DataStream

```

d. 希望能够按日期分目录存储

为了支持hive的分区处理，hdfs sink在将日志写入到hdfs的过程中，希望按照日期分目录存储。

```

1  a1.sinks.k1.hdfs.path = hdfs://hadoop01:9000/flux/reportTime=%Y-%m-%d

```

并且通过拦截器在日志头中增加timestamp头

```

1  a1.sources.r1.interceptors = i1
2  a1.sources.r1.interceptors.i1.type = regex_extractor
3  a1.sources.r1.interceptors.i1.serializers = s1
4  a1.sources.r1.interceptors.i1.serializers.s1.name = timestamp
5

```

e. 生成的目录时间不正确

配置hdfs采用的时区

```

1  a1.sinks.k1.hdfs.timeZone = GMT+8

```