

PRÁCTICA N°4.5

Actividad 4.5. Agente Flume con sumidero en HDFS

- Nombre y apellidos: Alvaro Lucio-Villegas De Cea



Índice

Vamos a preparar el Fichero de Configuración.	3
Prepara el fichero log.txt:	5
Crea el directorio de Volcado en HDFS:	6
Crea el Agente Flume:	6
Modifica el Fichero /temp/log.txt	8
Responde a las siguientes cuestiones:	9

Enunciado

Vamos a crear un agente Flume, de tal forma que registre eventos que vamos a generar modificando un el fichero de texto que estamos utilizando como si fuera un fichero de Log.

La idea es que se cree el fichero, y se vaya incrementando su contenido de forma dinámica, de tal forma que los cambios registrados en el fichero, no se muestren en la consola, sino que se vuelquen en ficheros de texto de HDFS. Como sabemos, el fichero de log es /tmp/log.txt.

Vamos a preparar el Fichero de Configuración.

- Utiliza el directorio flumeconf en la carpeta /home/cloudera, que ya utilizamos en la actividad 4.4
- Dentro de ese directorio, crea el fichero agentefichero.loghdfs.conf.
- Utiliza como nombre del agente agentefichero.loghdfs.
- Ese fichero debe contener todos los elementos de configuración que hemos visto en teoría. Vamos a concretarlos:

Como puede verse en el contenido del fichero que se muestra a continuación, aparece una configuración de propiedades del SINK, del sumidero mucho más elaborada, con cierto grado de complejidad. Este SINK o SUMIDERO es el destino en el que volcamos el contenido detectado por cada evento.

A continuación, una descripción de lo que representa cada elemento:

Propiedad	Valor	Significado
.sinks.sink1.type	hdfs	Tipo de destino, en este caso hdfs.
.sinks.sink1.hdfs.path	/flume/events/%y-%m-%d/%H%M	Ruta del cluster hdfs en la que depositará los eventos leídos de la fuente (source)
sink1.hdfs.filePrefix	events	Esto es un texto que utiliza flume para nombrar al fichero que genera en cada evento.
sinks.sink1.hdfs.round	true	Indica si debe redondear la marca de tiempo utilizada para construir los ficheros y carpetas que genera
sinks.sink1.hdfs.roundValue	1	Utiliza esta cantidad para partir las capturas de evento en ficheros y carpetas. La unidad que se utiliza es la que se indique en roundUnit , que en este caso son minutos.
sinks.sink1.hdfs.roundUnit	minute	Unidad para partir las carpetas
sinks.sink1.hdfs.useLocalTimeStamp	true	Usa el tiempo del servidor local, en lugar del tiempo que venga en la cabecera del evento.
sinks.sink1.hdfs.fileType	DataStream	Tipo de fichero generado. Se admiten los siguientes: SequenceFile, DataStream or CompressedStream

Fichero agenteficherologhdfs.conf:

```
# Definición de componentes del agente
agenteficherolog.sources = source1
agenteficherolog.sinks = sink1
agenteficherolog.channels = ch1

# Configuración de propiedades del source
agenteficherolog.sources.source1.type = exec
agenteficherolog.sources.source1.command = tail -F /tmp/log.txt
agenteficherolog.sources.source1.shell = /bin/bash -c

# Configuración de propiedades del sink
agenteficherologhdfs.sinks.sink1.type = hdfs
agenteficherologhdfs.sinks.sink1.hdfs.path = /flume/events/%y-%m-%d/%H%M
agenteficherologhdfs.sinks.sink1.hdfs.filePrefix = events
agenteficherologhdfs.sinks.sink1.hdfs.round = true
agenteficherologhdfs.sinks.sink1.hdfs.roundValue = 1
agenteficherologhdfs.sinks.sink1.hdfs.roundUnit = minute
agenteficherologhdfs.sinks.sink1.hdfs.useLocalTimeStamp = true
agenteficherologhdfs.sinks.sink1.hdfs.fileType = DataStream

# Configuración de un canal de tipo memoria
agenteficherolog.channels.ch1.type = memory
agenteficherolog.channels.ch1.capacity = 1000
agenteficherolog.channels.ch1.transactionCapacity = 100

# Vincular source y sink al canal creado
agenteficherolog.sources.source1.channels = ch1
agenteficherolog.sinks.sink1.channel = ch1
```

```
agentefichero loghdfs.conf  log.txt
# Definición de componentes del agente
agentefichero loghdfs.sources = source1
agentefichero loghdfs.sinks = sink1
agentefichero loghdfs.channels = ch1

# Configuración de propiedades del source
agentefichero loghdfs.sources.source1.type = exec
agentefichero loghdfs.sources.source1.command = tail -F /tmp/log.txt
agentefichero loghdfs.sources.source1.shell = /bin/bash -c

# Configuración de propiedades del sink
agentefichero loghdfs.sinks.sink1.type = hdfs
agentefichero loghdfs.sinks.sink1.hdfs.path = /flume/events/%y-%m-%d/%H%M
agentefichero loghdfs.sinks.sink1.hdfs.filePrefix = events
agentefichero loghdfs.sinks.sink1.hdfs.round = true
agentefichero loghdfs.sinks.sink1.hdfs.roundValue = 1
agentefichero loghdfs.sinks.sink1.hdfs.roundUnit = minute
agentefichero loghdfs.sinks.sink1.hdfs.useLocalTimeStamp = true
agentefichero loghdfs.sinks.sink1.hdfs.fileType = DataStream

# Configuración de un canal de tipo memoria
agentefichero loghdfs.channels.ch1.type = memory
agentefichero loghdfs.channels.ch1.capacity = 1000
agentefichero loghdfs.channels.ch1.transactionCapacity = 100

# Vincular source y sink al canal creado
agentefichero loghdfs.sources.source1.channels = ch1
agentefichero loghdfs.sinks.sink1.channel = ch1
```

Prepara el fichero log.txt:

- Crea el fichero log.txt en la ruta /tmp con el editor de texto de la máquina Cloudera.
- Añade un texto cualquiera
- Cierra el Editor

```
log.txt (/tmp) - gedit
File Edit View Search Tools Documents Help
Open Save Undo
agentefichero loghdfs.conf  log.txt
Alvaro Lucio
Esto es una prueba de Flume.
```

Crea el directorio de Volcado en HDFS:

- Crea el directorio /flume/events en HDFS, utilizando un terminal y las líneas de comando, mediante el comando “hdfs dfs.....”.

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /flume
[cloudera@quickstart ~]$ hdfs dfs -mkdir /flume/events
[cloudera@quickstart ~]$ hdfs dfs -ls
Found 2 items
drwxr-xr-x - cloudera cloudera          0 2023-03-20 13:00 _sqoop
drwxr-xr-x - cloudera cloudera          0 2023-03-20 12:36 dataset-Telefonia
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 10 items
drwxrwxrwx - hdfs      supergroup       0 2017-10-23 09:15 /benchmarks
drwxr-xr-x - cloudera supergroup       0 2023-04-10 12:09 /flume
drwxr-xr-x - hbase    supergroup       0 2023-03-13 12:16 /hbase
drwxr-xr-x - solr     solr             0 2017-10-23 09:18 /solr
drwxr-xr-x - cloudera supergroup       0 2023-03-13 12:24 /tablascolhdfs
drwxr-xr-x - cloudera supergroup       0 2023-03-29 12:02 /tablashdfs
drwxr-xr-x - root     supergroup       0 2023-03-29 11:39 /tablashdfsav
drwxrwxrwt - hdfs     supergroup       0 2023-02-22 12:34 /tmp
drwxr-xr-x - hdfs     supergroup       0 2017-10-23 09:17 /user
drwxr-xr-x - hdfs     supergroup       0 2017-10-23 09:17 /var
[cloudera@quickstart ~]$
```

Crea el Agente Flume:

Siguiendo las indicaciones que aparecen en teoría para crear el agente Flume, construye el

comando correspondiente, que tenga en cuenta lo siguiente:

- Nombre del Agente: agenteficherologhdfs
- Ruta del Directorio de Configuración: /home/cloudera/flumeconf
- Nombre del Fichero de configuración: agenteficherologhdfs.conf

```
flume-ng agent - agenteficherologhdfs -f
/home/cloudera/flumeconf/agenteficherologhdfs.conf
-Dflume.root.logger=INFO,console-Xmx512m
```

```

23/04/10 12:20:02 INFO node.PollingPropertiesFileConfigurationProvider: Configuration provider starting
23/04/10 12:20:02 INFO node.PollingPropertiesFileConfigurationProvider: Reloading configuration file:/home/cloudera/flumeconf/agentficerohloghdfs.conf
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Processing:sink1
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Added sinks: sink1 Agent: agentficerohloghdfs
23/04/10 12:20:02 INFO conf.FlumeConfiguration: Post-validation flume configuration contains configuration for agents: [agentficerohloghdfs]
23/04/10 12:20:02 INFO node.AbstractConfigurationProvider: Creating
23/04/10 12:20:02 INFO channel.DefaultChannelFactory: Creating instance of channel chl type memory
23/04/10 12:20:02 INFO node.AbstractConfigurationProvider: Created channel chl
23/04/10 12:20:02 INFO source.DefaultSourceFactory: Creating instance of source sourcecl, type exec
23/04/10 12:20:02 INFO sink.DefaultSinkFactory: Creating instance of sink sink1, type hdfs
23/04/10 12:20:02 INFO node.AbstractConfigurationProvider: Channel chl connected to [sourcecl, sink1]
23/04/10 12:20:02 INFO node.Application: Starting new configuration: {sourceRunners:[{source=EventDrivenSourceRunner; {source.org.apache.flume.source.MemoryCounterGroup=org.apache.flume.source.EventDrivenSourceRunner$MemoryCounterGroup; counterGroup={name:null; counters:{} } } ]}, channels:[{chl=org.apache.flume.channel.ChannelMemory}]}
23/04/10 12:20:02 INFO node.Application: Starting Channel chl
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: CHANNEL, name: chl: Successfully registered new MBean.
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Component type: CHANNEL, name: chl started
23/04/10 12:20:02 INFO node.Application: Starting Sink sink1
23/04/10 12:20:02 INFO node.Application: Starting Source sourcecl
23/04/10 12:20:02 INFO source.ExecSource: Exec source starting with command: tail -F /tmp/log.txt
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Monitor counter group Agr type: SINK, name: sink1: Successfully registered new MBean.
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Component type: SINK, name: sink1 started
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SOURCE, name: sourcecl: Successfully registered new MBean.
23/04/10 12:20:02 INFO instrumentation.MonitoredCounterGroup: Component type: SOURCE, name: sourcecl started
23/04/10 12:20:06 INFO hdfs.HDFSOutputStream: Serializer = TEXT useRawLocalFileSystem = false
23/04/10 12:20:06 INFO hdfs.BucketWriter: Creating /flume/events/23-04-10/1220/events.1681154406607.tmp
23/04/10 12:20:38 INFO hdfs.BucketWriter: Closing /flume/events/23-04-10/1220/events.1681154406607.tmp
23/04/10 12:20:38 INFO hdfs.HDFSOutputStream: Writing /flume/events/23-04-10/1220/events.1681154406607.tmp to /flume/events/23-04-10/1220/events.1681154406607
23/04/10 12:20:38 INFO hdfs.HDFSEventSink: Write callbac called.

```

A partir de este momento, el agente está escuchando en las operaciones que se hagan sobre el fichero /tmp/log.txt, y reaccionará a cada cambio del fichero de forma periódica, creando, si no existe una carpeta con el día actual, una carpeta con el minuto actual y dentro generará ficheros que tendrán nombres como "events.12331242142".

7

Modifica el Fichero /temp/log.txt

- Para modificar el fichero, vamos a ejecutar un código que añada líneas de forma periódica al fichero. Para eso, utilizaremos la siguiente instrucción, que añade una nueva línea con fecha y hora cada 5 segundos:

```
while true; do date >> /tmp/log.txt; sleep 5; done
```

```
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ while true;do date >>/tmp/log.txt;sleep 5;done
```

Comprobamos en HDFS que se están registrando todos los cambios con fecha y hora

```
[cloudera@quickstart ~]$ hdfs dfs -ls /flume/events/23-04-10/
Found 2 items
drwxr-xr-x - cloudera supergroup          0 2023-04-10 12:20 /flume/events/23-04-10/1220
drwxr-xr-x - cloudera supergroup          0 2023-04-10 12:25 /flume/events/23-04-10/1225
[cloudera@quickstart ~]$ hdfs dfs -ls /flume/events/23-04-10/1225
Found 2 items
-rw-r--r-- 1 cloudera supergroup      203 2023-04-10 12:25 /flume/events/23-04-10/1225/events.1681154705892
-rw-r--r-- 1 cloudera supergroup      29 2023-04-10 12:25 /flume/events/23-04-10/1225/events.1681154736018.tmp
[cloudera@quickstart ~]$
```

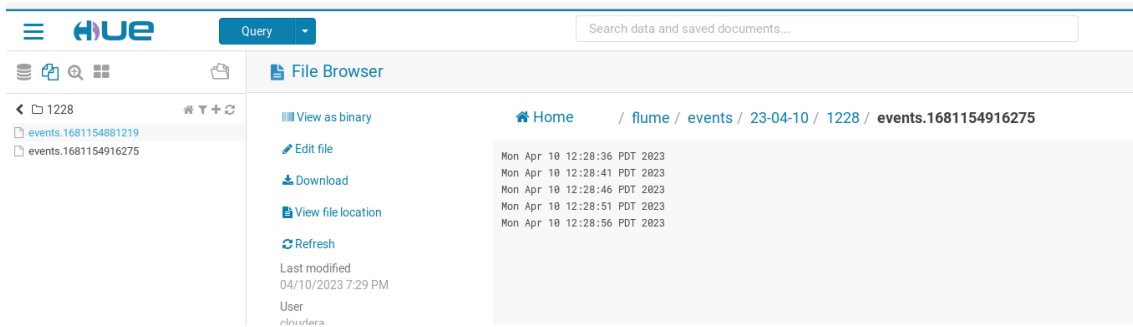

Responde a las siguientes cuestiones:

- Observa la terminal donde has lanzado el agente flume. ¿Qué ocurre?

-Se está mostrando toda la información sobre los sucesos del fichero.

- Observa el directorio HDFS desde el navegador, mediante la herramienta HUE.

Consulta el contenido de los directorios y de los ficheros.



- ¿Cada cuanto se genera un nuevo directorio en el directorio /flume/events del HDFS?

-Cada día

¿Y en el directorio /flume/events/%y-%m-%d?

-Cada Minuto

- ¿Cada cuanto se genera un nuevo fichero de evento? ¿A qué crees que se debe?

Revisa la propiedad rollInterval en la documentación de Flume

-Por defecto está puesto cada 30 segundos.Y esto sucede por el parámetro "rollinterval" (<https://flume.apache.org/FlumeUserGuide.html#flume-sinks>)

Cambia para que se generen 4 ficheros de eventos cada minuto.

Se agrega en el fichero de configuración el parámetro “rollintervall”

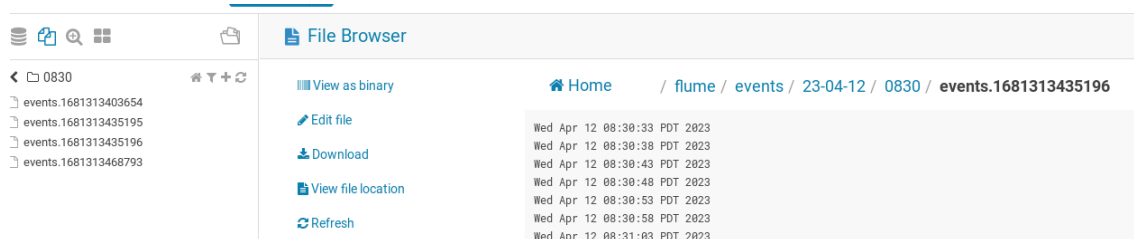
```
agentefichero loghdfs.conf x log.txt x
# Definición de componentes del agente
agentefichero loghdfs.sources = source1
agentefichero loghdfs.sinks = sink1
agentefichero loghdfs.channels = ch1

# Configuración de propiedades del source
agentefichero loghdfs.sources.source1.type = exec
agentefichero loghdfs.sources.source1.command = tail -F /tmp/log.txt
agentefichero loghdfs.sources.source1.shell = /bin/bash -c

# Configuración de propiedades del sink
agentefichero loghdfs.sinks.sink1.type = hdfs
agentefichero loghdfs.sinks.sink1.hdfs.path = /flume/events/%y-%m-%d/%H%M
agentefichero loghdfs.sinks.sink1.hdfs.filePrefix = events
agentefichero loghdfs.sinks.sink1.hdfs.round = true
agentefichero loghdfs.sinks.sink1.hdfs.roundValue = 1
agentefichero loghdfs.sinks.sink1.hdfs.roundUnit = minute
agentefichero loghdfs.sinks.sink1.hdfs.useLocalTimeStamp = true
agentefichero loghdfs.sinks.sink1.hdfs.fileType = DataStream
agentefichero loghdfs.sinks.sink1.hdfs.rollintervall = 14

# Configuración de un canal de tipo memoria
agentefichero loghdfs.channels.ch1.type = memory
agentefichero loghdfs.channels.ch1.capacity = 1000
agentefichero loghdfs.channels.ch1.transactionCapacity = 100

# Vincular source y sink al canal creado
agentefichero loghdfs.sources.source1.channels = ch1
agentefichero loghdfs.sinks.sink1.channel = ch1
```



- Cambia a 2 minutos la agrupación de eventos en carpetas. Prueba a hacerlo sin parar el agente flume ¿Qué ocurre?

Se almacenan más eventos en una sola carpeta.

```
agentfichero.loghdfs.conf
agentfichero.loghdfs.sources.source1.shell = /bin/bash -c

# Configuración de propiedades del sink
agentfichero.loghdfs.sinks.sink1.type = hdfs
agentfichero.loghdfs.sinks.sink1.hdfs.path = /flume/events/%y-%m-%d/%H%M
agentfichero.loghdfs.sinks.sink1.hdfs.filePrefix = events
agentfichero.loghdfs.sinks.sink1.hdfs.round = true
agentfichero.loghdfs.sinks.sink1.hdfs.roundValue = 2
agentfichero.loghdfs.sinks.sink1.hdfs.roundUnit = minute
agentfichero.loghdfs.sinks.sink1.hdfs.useLocalTimeStamp = true
agentfichero.loghdfs.sinks.sink1.hdfs.fileType = DataStream
agentfichero.loghdfs.sinks.sink1.hdfs.rollinterval = 4

# Configuración de un canal de tipo memoria
agentfichero.loghdfs.channels.ch1.type = memory
agentfichero.loghdfs.channels.ch1.capacity = 1000
agentfichero.loghdfs.channels.ch1.transactionCapacity = 100

# Vincular source y sink al canal creado
agentfichero.loghdfs.sources.source1.channels = ch1
agentfichero.loghdfs.sinks.sink1.channel = ch1
```

0827
0828
0829
0830
0832

- Modifica el proceso que añade líneas al fichero de log, para que escriba cada 10 segundos, en lugar de cada 5 segundos ¿Qué ocurre con los ficheros generados?

```
cloudera@quickstart ~]$ while true;do date >>/tmp/log.txt;sleep 10;done
```

Se generan menos entradas en el log.

0834
events.1681313648475

View as binary
Edit file
Download
View file location
Refresh

Home / flume / events / 23-04-12 / 0834 / events.1681313648475

Wed Apr 12 08:34:08 PDT 2023
Wed Apr 12 08:34:18 PDT 2023
Wed Apr 12 08:34:28 PDT 2023
Wed Apr 12 08:34:38 PDT 2023

-Cámbialo a 2 segundos y explica lo que ocurre con los ficheros generados en esta ocasión

```
[cloudera@quickstart ~]$ while true;do date >>/tmp/log.txt;sleep 2;done
```

Se generan más entradas en el log.

The screenshot shows the Cloudera interface for a file named `events.1681319405538.tmp`. The sidebar on the left lists the file and its parent directory `1010`. The central panel provides actions for the file: `View as binary`, `Edit file`, `Download`, `View file location`, and `Refresh`. It also shows the file's last modified date as `04/12/2023 5:10 PM` and the user `cloudera`. The right panel displays the file's content as a list of timestamps: `Wed Apr 12 10:10:23 PDT 2023` through `Wed Apr 12 10:10:42 PDT 2023`.