

实时数据生成器

Python Socket 实时数据生成器

以下是一个使用Python编写的简单Socket服务器，可以生成实时数据并发送到指定的端口（如10050），供Flume采集：

```

1  import socket
2  import time
3  import random
4
5  # 创建Socket服务器
6  def start_data_generator():
7      host = '0.0.0.0' # 接收所有可用IP地址
8      port = 10050 # 端口号
9
10     # 创建socket对象, 指定协议和数据流类型
11     server_socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
12     server_socket.bind((host, port))
13
14     # 允许的最大连接数
15     server_socket.listen(5)
16
17     print(f"实时数据生成器已启动, 正在监听 {port} 端口...")
18
19     while True:
20         client_socket, addr = server_socket.accept()
21         print(f"连接来自: {addr}")
22
23         try:
24             while True:
25                 # 生成模拟订单数据
26                 order_id = random.randint(1000, 9999)
27                 product_id = random.randint(1, 100)
28                 quantity = random.randint(1, 10)
29                 price = round(random.uniform(10.0, 100.0), 2)
30                 timestamp = time.strftime('%Y-%m-%d %H:%M:%S', time.localtime())
31
32                 # 生成数据格式, 例如: order_id, product_id, quantity, price,
33                 # timestamp
34                 data = f"{order_id},{product_id},{quantity},{price},{timestamp}\n"
35
36                 # 向客户端发送数据
37                 client_socket.sendall(data.encode('utf-8'))
38                 print(f"发送数据: {data.strip()}")
39
40                 # 模拟延迟
41                 time.sleep(1)
42
43         except BrokenPipeError:
44             print(f"客户端 {addr} 已断开连接。")

```

```
44         finally:
45             client_socket.close()
46
47     # 启动数据生成器
48     if __name__ == "__main__":
49         start_data_generator()
```

说明：

1. 该脚本将启动一个Socket服务器，监听端口10050。
2. 每秒生成一条模拟订单数据，格式为 `order_id, product_id, quantity, price, times tamp`，并发送给连接的客户端。
3. 运行此脚本后，Flume可以通过配置一个 `NetcatSource` 来连接并接收数据。

可以根据需要调整数据生成的格式和字段。

Scala版本的实时数据生成器

以下是一个使用Scala编写的简单Socket服务器，可以生成实时数据并发送到指定的端口（如10050），供Flume采集：

```

1  import java.io.PrintStream
2  import java.net.ServerSocket
3  import scala.util.Random
4  import java.time.LocalDateTime
5  import java.time.format.DateTimeFormatter
6  import scala.concurrent.{Future, ExecutionContext}
7  import scala.concurrent.ExecutionContext.Implicits.global
8
9  object RealTimeDataGenerator {
10
11  def main(args: Array[String]): Unit = {
12      val server = new ServerSocket(10050)
13      println("Scala 实时数据生成器已启动, 监听端口 10050...")
14
15      while (true) {
16          val socket = server.accept()
17          println(s"客户端连接: ${socket.getInetAddress}")
18
19          // 使用Future处理并发连接
20          Future {
21              val out = new PrintStream(socket.getOutputStream)
22
23              try {
24                  while (true) {
25                      // 生成随机订单数据
26                      val orderId = Random.nextInt(10000)
27                      val productId = Random.nextInt(100)
28                      val quantity = Random.nextInt(10) + 1
29                      val price = Random.nextDouble() * 100
30                      val timestamp = LocalDateTime.now.format(DateTimeFormatter.ofP
attern("yyyy-MM-dd HH:mm:ss"))
31
32                      // 拼接数据
33                      val data = s"$orderId,$productId,$quantity,$price,$timestamp
\n"
34
35                      out.print(data)
36                      out.flush()
37
38                      println(s"发送数据: $data")
39
40                      // 模拟数据生成的延迟
41                      Thread.sleep(1000)
42                  }
43              } catch {
44                  case e: Exception => println(s"连接出错: ${e.getMessage}")
45              } finally {

```

```
45         out.close()
46         socket.close()
47     }
48 }
49 }
50 }
51 }
```

说明：

1. 该程序会在本地 `10050` 端口启动一个Socket服务，客户端连接后，生成并发送类似 `order_id,product_id,quantity,price,timestamp` 的订单数据。
2. 每秒生成一条数据，模拟订单系统的实时数据流。
3. 数据格式为： `订单ID，产品ID，购买数量，价格，时间戳` 。

Flume配置文件示例

假设使用Flume来采集上述Socket数据并将其发送到Kafka中的一个名为 `order` 的Topic，以下是Flume的配置文件示例：

flume.conf：

```

1  # Define the agent name
2  agent.sources = source1
3  agent.sinks = sink1
4  agent.channels = channel1
5
6  # Source configuration (using Netcat to receive socket data from the Scala generator)
7  agent.sources.source1.type = netcat
8  agent.sources.source1.bind = 0.0.0.0
9  agent.sources.source1.port = 10050
10
11 # Sink configuration (send data to Kafka)
12 agent.sinks.sink1.type = org.apache.flume.sink.kafka.KafkaSink
13 agent.sinks.sink1.kafka.bootstrap.servers = localhost:9092
14 agent.sinks.sink1.kafka.topic = order
15 agent.sinks.sink1.kafka.partitioner.class = org.apache.kafka.clients.producer.internals.DefaultPartitioner
16 agent.sinks.sink1.kafka.producer.acks = 1
17
18 # Channel configuration (using memory channel)
19 agent.channels.channel1.type = memory
20 agent.channels.channel1.capacity = 10000
21 agent.channels.channel1.transactionCapacity = 100
22
23 # Bind source and sink to channel
24 agent.sources.source1.channels = channel1
25 agent.sinks.sink1.channel = channel1

```

说明：

1. **Source部分**：使用 `NetcatSource` 来监听 `10050` 端口，从Scala生成器中接收数据。
2. **Sink部分**：使用 `KafkaSink` 将数据发送到Kafka中的 `order` 主题，Kafka服务器运行在 `localhost:9092`，你可以根据需要调整服务器地址和端口。
3. **Channel部分**：使用 `memory channel` 作为Flume数据流动的桥梁，确保数据从Source流向Sink。

启动Flume：

确保Flume已经安装并配置好Kafka连接，运行以下命令启动Flume：

```

1  flume-ng agent --conf ./conf --conf-file flume.conf --name agent -Dflume.root.logger=INFO,console

```

验证：

1. 数据发送到Kafka：使用Kafka自带的命令行工具来查看数据：

```
1 kafka-console-consumer --bootstrap-server localhost:9092 --topic order --from-beginning
```

2. 多路复用备份到HDFS：还可以配置Flume的多路复用，将数据同时写入Kafka和HDFS。

为了配置Flume将数据同时写入Kafka和HDFS，可以使用Flume的多路复用功能。通过多路复用，Flume可以将相同的数据同时发送到多个Sink，如Kafka和HDFS。

修改后的 Flume 配置文件示例：

flume-multiplex.conf：

```

1  # 定义 agent 名称
2  agent.sources = source1
3  agent.sinks = sink1 sink2
4  agent.channels = channel1 channel2
5
6  # Source 配置 (Netcat 源接收来自 Scala 生成器的 Socket 数据)
7  agent.sources.source1.type = netcat # 使用 netcat 作为 source 类型
8  agent.sources.source1.bind = 0.0.0.0 # 绑定到所有可用的网络接口
9  agent.sources.source1.port = 10050 # 监听端口 10050 接收数据
10
11 # Channel 选择
12 agent.sources.source1.channelSelector.type = replicating # 使用复制模式, 数
    据同时发送到多个通道
13 agent.sources.source1.channels = channel1 channel2 # 绑定到两个通道 (channe
    l1 和 channel2)
14
15 # Sink 1 配置 (将数据发送到 Kafka)
16 agent.sinks.sink1.type = org.apache.flume.sink.kafka.KafkaSink # 使用 Kaf
    kaSink 作为输出目的地
17 agent.sinks.sink1.kafka.bootstrap.servers = localhost:9092 # Kafka 集群的
    bootstrap 服务器地址
18 agent.sinks.sink1.kafka.topic = order # Kafka 主题名称为 order
19 agent.sinks.sink1.kafka.partitioner.class = org.apache.kafka.clients.produ
    cer.internals.DefaultPartitioner # 使用默认的 Kafka 分区器
20 agent.sinks.sink1.kafka.producer.acks = 1 # 确认机制, 1 表示等待 Kafka 服务器
    的确认
21
22 # Sink 2 配置 (将数据备份到 HDFS)
23 agent.sinks.sink2.type = hdfs # 使用 HDFS 作为输出目的地
24 agent.sinks.sink2.hdfs.path = hdfs://namenode/user/test/flumebackup/ # HD
    FS 备份路径
25 agent.sinks.sink2.hdfs.fileType = DataStream # 数据文件类型为流式数据
26 agent.sinks.sink2.hdfs.writeFormat = Text # 写入数据格式为文本格式
27 agent.sinks.sink2.hdfs.batchSize = 1000 # 批处理大小, 1000 条事件后进行写入
28 agent.sinks.sink2.hdfs.rollInterval = 30 # 每隔 30 秒滚动生成一个新的文件
29 agent.sinks.sink2.hdfs.rollSize = 0 # 文件大小达到指定值后滚动 (此处为 0, 表示
    忽略)
30 agent.sinks.sink2.hdfs.rollCount = 0 # 写入指定数量的事件后滚动 (此处为 0, 表示
    忽略)
31
32 # 通道 1 配置 (用于 Kafka Sink)
33 agent.channels.channel1.type = memory # 使用内存通道
34 agent.channels.channel1.capacity = 10000 # 通道容量为 10000 条事件
35 agent.channels.channel1.transactionCapacity = 100 # 单次事务容量为 100 条事
    件
36

```



```

37 # 通道 2 配置 (用于 HDFS Sink)
38 agent.channels.channel2.type = memory # 使用内存通道
39 agent.channels.channel2.capacity = 10000 # 通道容量为 10000 条事件
40 agent.channels.channel2.transactionCapacity = 100 # 单次事务容量为 100 条事件
41
42 # 将 source 和 sink 绑定到各自的通道
43 agent.sources.source1.channels = channel1 channel2 # 将 source 绑定到两个通道
44 agent.sinks.sink1.channel = channel1 # 将 Kafka Sink 绑定到 channel1
45 agent.sinks.sink2.channel = channel2 # 将 HDFS Sink 绑定到 channel2
46

```

配置说明：

1. Source 部分：

- 使用 `NetcatSource` 来监听 `10050` 端口，接收来自Scala生成器的Socket数据。

2. ChannelSelector：

- 使用 `replicating` 类型的选择器，表示Flume会将数据同时发送到两个Channel（`channel1` 和 `channel2`），从而实现多路复用。

3. Sink 1 (KafkaSink)：

- 将数据写入Kafka的 `order` 主题。Kafka运行在 `localhost:9092`。

4. Sink 2 (HDFS Sink)：

- 将数据备份到HDFS路径 `hdfs://namenode/user/test/flumebackup/`。
- `rollInterval=30` 表示每隔30秒会滚动生成一个新的文件。
- `rollSize` 和 `rollCount` 设置为 `0`，表示根据时间滚动生成文件，而不是基于文件大小或事件数量。

5. Channels：

- 为Kafka Sink和HDFS Sink分别配置了两个独立的Memory Channel。

启动 Flume：

```

1 flume-ng agent --conf ./conf --conf-file flume-multiplex.conf --name agent
  -Dflume.root.logger=INFO,console

```

验证配置：

1. 验证 Kafka 数据：

使用Kafka的命令行工具来验证数据是否成功发送到 `order` 主题：

```
1 kafka-console-consumer --bootstrap-server localhost:9092 --topic order --from-beginning
```

2. 验证 HDFS 备份：

可以在HDFS中查看Flume备份的文件，确认数据是否备份成功：

```
1 hdfs dfs -ls /user/test/flumebackup/  
2 hdfs dfs -cat /user/test/flumebackup/<file_name>
```

通过以上配置，Flume将会将数据同时写入Kafka和HDFS。