

UD 7 - Apache Kafka y Spark

1. Introducción

Para entender y aprender como podemos usar Kafka y Spark en un sistema de Big Data, lo vamos a realizar mediante un ejemplo

2. Ejemplo 6. Kafka y Spark

Imagina que vas a trabajar para un gran equipo de fútbol en el análisis de datos para obtener business intelligence para que durante los partidos, podamos servir y dar soporte al stack técnico.

Vamos a realizar un pequeño ejemplo de concepto, donde vamos a generar datos sintéticos de las posiciones de los jugadores durante el partido, para seguidamente obtener información como: distancia recorrida, velocidades, sprints, etc; en tiempo real, para que ayude al stack técnico en su toma de decisiones.

Vamos a establecer un cluster de Kafka para recoger los eventos en tiempo real, que dispondrá los datos para que Spark los recoja en tiempo real en un dataset y pueda aplicarle Pandas para obtener estos datos.

✓ Configuración cluster Kafka

Vamos a **configurar el cluster de Kafka usando mode KRaft** teniendo en cuenta que estamos en un **entorno de pruebas y no de producción**, para este ejemplo de concepto, vamos a configurar un **cluster con 3 servidores (en un sólo nodo o máquina) con 1 controller y 2 brokers**

2.1 Requisitos

Vamos a configurar un cluster las siguientes características:

- 1 topic jugadores con 2 particiones
- 1 factor de replica de 2
- 1 nodo con 2 brokers y 1 controller

2.2 Configuración del Clúster de Kafka

- Consideraciones previas:

1. Vamos a establecer todos los archivos de configuración en una carpeta de ejemplo llamada ejemplo6, que en mi caso estará alojada en `/opt/kafka/ejemplo6`
2. Dado que todas las instancias se ejecutarán en el mismo nodo, es crucial asignar puertos únicos y directorios de log para cada broker y el controller.

- Configuración:

1. Para el controller, debemos usar como base la configuración de propiedades de controller de kafka (KRaft mode) que se encuentran `config/kraft/controller.properties`
2. Para cada broker, necesitaremos crear un archivo de configuración por separado. Para ello debemos usar como base la configuración de propiedades de brokers de kafka que se encuentran `config/kraft/broker.properties`
3. Creamos los directorios necesarios para nuestro `ejemplo6`

```
1 mkdir -p /opt/kafka/ejemplo6/config
2 mkdir -p /opt/kafka/ejemplo6/logs
```

4. Hacemos 2 y 1 copia de los ficheros correspondientes de configuración para cada uno

```
1 cp config/kraft/controller.properties /opt/kafka/ejemplo6/config/controller1.properties
2 cp config/kraft/broker.properties /opt/kafka/ejemplo6/config/broker1.properties
3 cp config/kraft/broker.properties /opt/kafka/ejemplo6/config/broker2.properties
```

5. Asignamos la configuración al controller

controller1.properties

```

1 # Server Basics
2 process.roles=controller
3 node.id=1
4 controller.quorum.voters=1@localhost:9093
5 # Socket Server Settings
6 listeners=CONTROLLER://localhost:9093
7 controller.listener.names=CONTROLLER
8 # Log Basics
9 log.dirs=/opt/kafka/ejemplo6/logs/controller1

```

6. Asignamos la siguiente configuración para cada broker

✓ Spark como consumer (nodo master de Spark)

Hay que tener en cuenta:

- Vamos a conectarnos con un cliente(consumer) que va a consumir datos en streaming con SPARK
- Por seguridad, kafka no permite esa conexión del cliente(consumer) desde localhost
- Teniendo en cuenta nuestra configuración del cluster (ip del master de Spark es 192.168.11.10), debemos configurarla correctamente en los brokers para que tengan acceso.
- La conexión se hace a través del master de spark, y la ejecución de la aplicación en Spark se reparte entre los workers del cluster, aunque estén en diferentes ips

broker1

broker1.properties

```

1 # Server Basics
2 process.roles=broker
3 node.id=2
4 controller.quorum.voters=1@localhost:9093
5 # Socket Server Settings
6 listeners=PLAINTEXT://192.168.11.10:9094
7 inter.broker.listener.name=PLAINTEXT
8 advertised.listeners=PLAINTEXT://192.168.11.10:9094
9 listener.security.protocol.map=CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,SSL:SSL,SASL_PLAINTEXT:SASL_PLAINTEXT,SASL_SSL:SASL_SSL
10 # Log Basics
11 log.dirs=/opt/kafka/ejemplo6/logs/broker1

```

broker2

broker2.properties

```

1 # Server Basics
2 process.roles=broker
3 node.id=3
4 controller.quorum.voters=1@localhost:9093
5 # Socket Server Settings
6 listeners=PLAINTEXT://192.168.11.10:9095
7 inter.broker.listener.name=PLAINTEXT
8 advertised.listeners=PLAINTEXT://192.168.11.10:9095
9 # Log Basics
10 log.dirs=/opt/kafka/ejemplo6/logs/broker2

```

⚠ Configuración con varios consumidores

En el caso de que tuviéramos que configurar más de un cliente en diferentes máquinas, tendríamos que cambiar la configuración de nuestros brokers. Recordando para que sirven las siguiente propiedades:

- `listeners`: son las interfaces a las que se conecta Kafka.
- `advertised.listeners`: es la forma en que los clientes pueden conectarse.
- `listener.security.protocol.map`: define pares clave/valor para el protocolo de seguridad a utilizar por nombre de listener.

1) Opción 1

Abrir el acceso en `listeners` configurando la ip como `0.0.0.0:puerto`. Esto no sería la mejor solución, ya que debemos filtrar el acceso desde `advertised.listeners` y este valor, por configuración propia de Kafka, **debe ser un subconjunto de `listeners`**

2) Opción 2

Indicamos por ip cada uno de los interfaces que damos acceso (**solo un interface por puerto**) a los brokers de Kafka (`listeners`), y filtramos (`advertised.listeners`). Pero no hay que olvidar indicar el protocolo por el cual te conectas, que hay que indicarlo en `listener.security.protocol.map`, que también debe ser único. Imaginando que tenemos varios clientes de acceso, la configuración debería ser en ese caso

broker1

broker1.properties

```
1 # Server Basics
2 process.roles=broker
3 node.id=2
4 controller.quorum.voters=1@localhost:9093
5 # Socket Server Settings
6 listeners=BROKER://localhost:9094,LISTENER_CLIENTE1://nodo1:11094,LISTENER_CLIENTE2://nodo2:12094,LISTENER_CLIENTE3://nodo3:13094
7 inter.broker.listener.name=BROKER
8 advertised.listeners=BROKER://localhost:9094,LISTENER_CLIENTE1://nodo1:11094,LISTENER_CLIENTE2://nodo2:12094,LISTENER_CLIENTE3://nodo3:13094
9 listener.security.protocol.map=CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,SSL:SSL,SASL_PLAINTEXT:SASL_PLAINTEXT,SASL_SSL:SASL_SSL,
NTEXT
10 # Log Basics
11 log.dirs=/opt/kafka/ejemplo6/logs/broker1
```

broker2

broker2.properties

```
1 # Server Basics
2 process.roles=broker
3 node.id=3
4 controller.quorum.voters=1@localhost:9093
5 # Socket Server Settings
6 listeners=BROKER://localhost:9095,LISTENER_CLIENTE1://nodo1:11095,LISTENER_CLIENTE2://nodo2:12095,LISTENER_CLIENTE3://nodo3:13095
7 inter.broker.listener.name=BROKER
8 advertised.listeners=BROKER://localhost:9095,LISTENER_CLIENTE1://nodo1:11095,LISTENER_CLIENTE2://nodo2:12095,LISTENER_CLIENTE3://nodo3:13095
9 listener.security.protocol.map=CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,SSL:SSL,SASL_PLAINTEXT:SASL_PLAINTEXT,SASL_SSL:SASL_SSL,
NTEXT
10 # Log Basics
11 log.dirs=/opt/kafka/ejemplo6/logs/broker2
```

7. Iniciamos Kafka con KRaft

```
1 #Genera un cluster UUID
2 KAFKA_CLUSTER_ID="$(bin/kafka-storage.sh random-uuid)"
3 echo $KAFKA_CLUSTER_ID
4
5 #Formateamos los directorios de log
6 bin/kafka-storage.sh format -t $KAFKA_CLUSTER_ID -c /opt/kafka/ejemplo6/config/controller1.properties
7 bin/kafka-storage.sh format -t $KAFKA_CLUSTER_ID -c /opt/kafka/ejemplo6/config/broker1.properties
8 bin/kafka-storage.sh format -t $KAFKA_CLUSTER_ID -c /opt/kafka/ejemplo6/config/broker2.properties
```

8. Iniciamos los server(1 controller y 2 brokers) cada uno en una terminal

```
1 #Ejecuta el servidor Kafka
2 bin/kafka-server-start.sh /opt/kafka/ejemplo6/config/controller1.properties
3 bin/kafka-server-start.sh /opt/kafka/ejemplo6/config/broker1.properties
4 bin/kafka-server-start.sh /opt/kafka/ejemplo6/config/broker2.properties
```

2.3 Creación del Topic

1. Creamos el topic con factor de replica 2 y 2 particiones. El topic debe conectarse a un broker.

```
1 bin/kafka-topics.sh --create --topic player-position --bootstrap-server 192.168.11.10:9094 --replication-factor 2 --
partitions 2
```

2. Podemos ver la descripción del topic creado

```
1 bin/kafka-topics.sh --describe --topic player-position --bootstrap-server 192.168.11.10:9094
```

```
1 Topic: player-position TopicId: DDEeIaCeQMaewIuwSBNQtW PartitionCount: 2 ReplicationFactor: 2 Configs:
segment.bytes=1073741824
2 Topic: player-position Partition: 0 Leader: 3 Replicas: 3,2 Isr: 3,2 Elr: LastKnownElr:
3 Topic: player-position Partition: 1 Leader: 2 Replicas: 2,3 Isr: 2,3 Elr: LastKnownElr:
```

2.4 Productor

1. Creamos el programa que va a generar los datos sintéticos.
2. Vamos a intentar simular datos lo más reales posibles a los movimientos de un jugador de futbol.
3. Suponemos que el campo es de 100x65 metros
4. El punto (0,0) será la esquina inferior izquierda y el punto (100,65) será la esquina superior derecha.
5. Generaremos 10 eventos por segundo

generate_synthetic_data.py

```
1 import random
2 import time
3 from datetime import datetime, timedelta
4 import math
5
6 # Parámetros del campo y jugadores
7 FIELD_LENGTH = 100
8 FIELD_WIDTH = 65
9 NUM_PLAYERS = 11
10 TIMESTAMPS_PER_SECOND = 10
11 DURATION_MINUTES = 45
12
13 # Velocidades en metros por segundo
14 WALK_SPEED = 1.4
15 RUN_SPEED = 5.0
16 SPRINT_SPEED = 8.0
17
18 def generate_initial_positions():
19     """Genera posiciones iniciales razonables para los jugadores."""
20     positions = {}
21     for player_id in range(1, NUM_PLAYERS + 1):
22         if player_id == 1:
23             # Portero
24             positions[player_id] = (random.uniform(0, 5), random.uniform(FIELD_WIDTH / 2 - 5, FIELD_WIDTH / 2 + 5))
25         else:
26             # Jugadores de campo
27             positions[player_id] = (random.uniform(0, FIELD_LENGTH), random.uniform(0, FIELD_WIDTH))
28     return positions
29
30 def move_player(position, speed):
31     """Mueve al jugador a una nueva posición basada en la velocidad y dirección."""
32     angle = random.uniform(0, 360) # Dirección aleatoria
33     distance = speed / TIMESTAMPS_PER_SECOND
34     delta_x = distance * math.cos(math.radians(angle))
35     delta_y = distance * math.sin(math.radians(angle))
36
37     new_x = max(0, min(FIELD_LENGTH, position[0] + delta_x))
38     new_y = max(0, min(FIELD_WIDTH, position[1] + delta_y))
39
40     return (new_x, new_y)
41
42 def generate_player_position_data(match_id, player_id, position):
43     timestamp = datetime.utcnow().isoformat()
44     return {
45         "match_id": match_id,
46         "player_id": player_id,
47         "timestamp": timestamp,
48         "position_x": position[0],
49         "position_y": position[1]
50     }
51
52 def simulate_player_movement(positions):
53     """Simula el movimiento de los jugadores en el campo."""
54     new_positions = {}
55     for player_id, position in positions.items():
56         if player_id == 1:
```

```

57         # Portero: se mueve dentro del área de gol
58         new_positions[player_id] = move_player(position, random.uniform(0, WALK_SPEED))
59     else:
60         # Jugadores de campo: pueden caminar, correr o sprintar
61         speed = random.choice([WALK_SPEED, RUN_SPEED, SPRINT_SPEED])
62         new_positions[player_id] = move_player(position, speed)
63     return new_positions
64
65 def main():
66     positions = generate_initial_positions()
67     end_time = datetime.utcnow() + timedelta(minutes=DURATION_MINUTES)
68
69     while datetime.utcnow() < end_time:
70         for _ in range(TIMESTAMPS_PER_SECOND):
71             for player_id in range(1, NUM_PLAYERS + 1):
72                 data = generate_player_position_data("match_1", player_id, positions[player_id])
73                 print(data) # Aquí enviaríamos los datos a Kafka en lugar de imprimirlos
74                 positions = simulate_player_movement(positions)
75                 time.sleep(1 / TIMESTAMPS_PER_SECOND)
76
77 if __name__ == "__main__":
78     main()

```

2. Creamos el productor con `KafkaProducer`

producer_positions.py

```

1  from kafka import KafkaProducer
2  import json
3  from generate_synthetic_data import generate_player_position_data, generate_initial_positions,
simulate_player_movement
4  import time
5  from datetime import datetime, timedelta
6
7  # Parámetros del campo y jugadores
8  NUM_PLAYERS = 11
9  TIMESTAMPS_PER_SECOND = 10
10 DURATION_MINUTES = 45
11
12 def main():
13     match_id = "match_" + datetime.utcnow().strftime("%Y%m%d%H%M%S")
14     producer = KafkaProducer(bootstrap_servers='192.168.11.10:9094',
15                             value_serializer=lambda v: json.dumps(v).encode('utf-8'))
16
17     positions = generate_initial_positions()
18     end_time = datetime.utcnow() + timedelta(minutes=DURATION_MINUTES)
19
20     while datetime.utcnow() < end_time:
21         for _ in range(TIMESTAMPS_PER_SECOND):
22             for player_id in range(1, NUM_PLAYERS + 1):
23                 data = generate_player_position_data(match_id, player_id, positions[player_id])
24                 producer.send('player-position', value=data)
25                 print(f"Sent: {data}")
26                 positions = simulate_player_movement(positions)
27                 time.sleep(1 / TIMESTAMPS_PER_SECOND)
28
29 if __name__ == "__main__":
30     main()

```

3. Aprovecharemos la *Consumer API* de Kafka para ver está consumiendo los datos correctamente una vez lanzada la aplicación

```

1  bin/kafka-console-consumer.sh --topic player-position --from-beginning --bootstrap-server 192.168.11.10:9094

```

2.5 Configuración del Clúster de Spark

Aprovecharemos las máquinas master y los 3 nodos que ya tenemos preparadas para lanzar el cluster de Spark. Consulta la documentación de Spark para cualquier duda

1. Código del programa

- a. Consumirá los eventos en tiempo real de kafka
- b. Lo almacena en un dataset
- c. Obtenemos ETL a través de Pandas
- d. Mostramos el resultado por consola:
 - i. El calculo de cada microBatch
 - ii. El acumulado del (último) partido.

2. Para una correcta configuración, nos basaremos en la documentación oficial de Apache Spark y su integración con Kafka: [Structured Streaming + Kafka Integration Guide](#)

Spark_PlayerPositionProcessing.py

```

1  from pyspark.sql import SparkSession
2  from pyspark.sql.functions import col, from_json
3  from pyspark.sql.types import StructType, StructField, IntegerType, StringType, FloatType
4  import pandas as pd
5  import numpy as np
6
7  # Definir el esquema de los datos
8  schema = StructType([
9      StructField("match_id", StringType(), True),
10     StructField("player_id", IntegerType(), True),
11     StructField("timestamp", StringType(), True),
12     StructField("position_x", FloatType(), True),
13     StructField("position_y", FloatType(), True)
14 ])
15
16 # Crear la sesión de Spark
17 spark = SparkSession.builder \
18     .appName("KafkaDataConsumer") \
19     .getOrCreate()
20
21 # Set Spark logging level to ERROR to avoid varios otros logs en consola.
22 spark.sparkContext.setLogLevel("ERROR")
23
24 # Leer los datos de Kafka
25 df = spark \
26     .readStream \
27     .format("kafka") \
28     .option("kafka.bootstrap.servers", "192.168.11.10:9094") \
29     .option("subscribe", "player-position") \
30     .option("startingOffsets", "earliest") \
31     .load()
32
33 player_position_df = df.selectExpr("CAST(value AS STRING)") \
34     .select(from_json(col("value"), schema).alias("data")) \
35     .select("data.*")
36
37 # Crear un DataFrame acumulativo vacío
38 accumulated_schema = StructType([
39     StructField("match_id", StringType(), True),
40     StructField("player_id", IntegerType(), True),
41     StructField("total_distance", FloatType(), True)
42 ])
43 accumulated_df = spark.createDataFrame([], accumulated_schema)
44
45 # Variable global para el 'match_id' actual
46 current_match_id = None
47
48 # Definir una función para procesar cada microbatch
49 def process_batch(df, epoch_id):
50     global accumulated_df, current_match_id
51
52     # Convertir a Pandas DataFrame
53     pd_df = df.toPandas()
54
55     # Obtener el 'match_id' del batch actual
56     if not pd_df.empty:
57         match_id = pd_df['match_id'].iloc[0]
58         current_match_id = match_id
59
60     # Calcular la distancia recorrida por cada jugador
61     pd_df['prev_position_x'] = pd_df.groupby('player_id')['position_x'].shift(1)
62     pd_df['prev_position_y'] = pd_df.groupby('player_id')['position_y'].shift(1)
63
64     pd_df['distance'] = np.sqrt((pd_df['position_x'] - pd_df['prev_position_x']) ** 2 +
65                               (pd_df['position_y'] - pd_df['prev_position_y']) ** 2)
66
67     pd_df['distance'] = pd_df['distance'].fillna(0)
68
69     # Sumar la distancia total por jugador
70     total_distance_df = pd_df.groupby(['match_id', 'player_id'])['distance'].sum().reset_index()
71
72     # Convertir de nuevo a DataFrame de Spark
73     spark_df = spark.createDataFrame(total_distance_df)
74
75     # Unir con el DataFrame acumulativo
76     accumulated_df = accumulated_df.union(spark_df).groupBy("match_id", "player_id").agg({"total_distance":
77 "sum"}).withColumnRenamed("sum(total_distance)", "total_distance")
78
79     # Filtrar el acumulado para el 'match_id' actual

```

```

79     current_match_df = accumulated_df.filter(col("match_id") == current_match_id)
80
81     # Mostrar el resultado del batch actual y el acumulado del último partido
82     print("Batch Result:")
83     spark_df.show()
84     print(f"Resultado acumulado para match_id = {current_match_id}:")
85     current_match_df.show()
86
87 # Escribir los datos procesados usando el método foreachBatch
88 query = player_position_df \
89     .writeStream \
90     .foreachBatch(process_batch) \
91     .outputMode("append") \
92     .start()
93
94 query.awaitTermination()

```

2. Lanzamos spark master y los workers del cluster.

```

1  hadoop@master:/opt/hadoop-3.4.1/spark-3.5.4$ ./sbin/start-master.sh
2  hadoop@master:/opt/hadoop-3.4.1/spark-3.5.4$ ./sbin/start-workers.sh

```

3. Recuerda que debes tener instalado correctamente pandas. Si no es el caso, hay que instalarlo en cada uno de los nodos

```

1  pip3 install pandas
2  pip3 install pyarrow

```

4. Lanzamos el programa. Siguiendo la [documentación oficial](#), añadimos el paquete del conector de Kafka

```

1  spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.4 --master spark://192.168.11.10:7077
/opt/kafka/ejemplo6/Spark_PlayerPositionProcessing.py

```

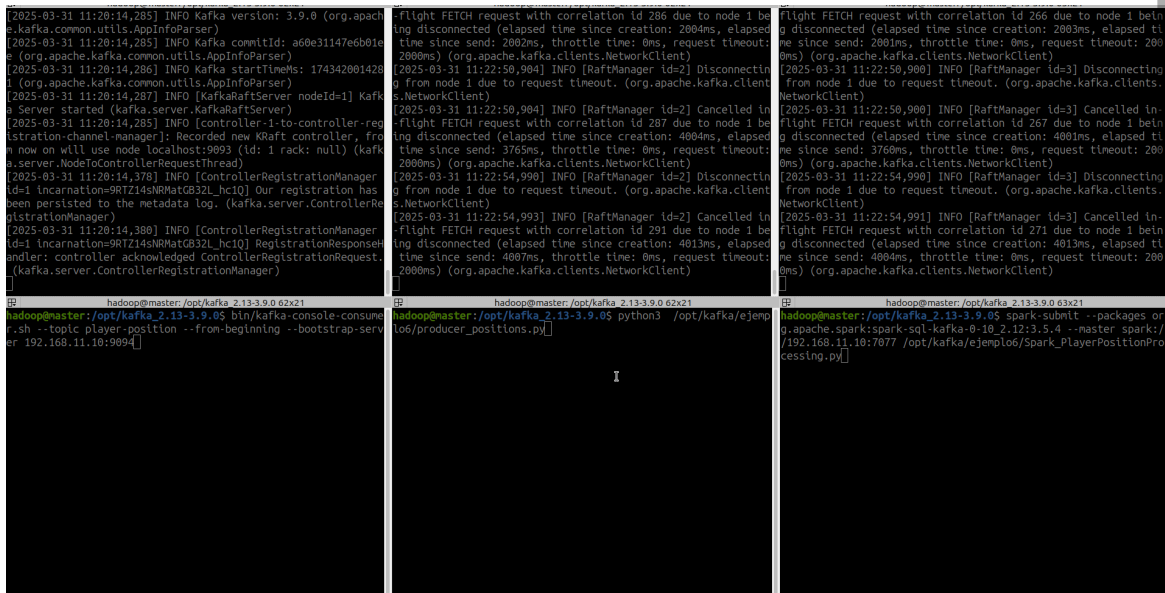
2.6 Ejecución de la aplicación

1. Lanzamos el productor

```

1  python3 /opt/kafka/ejemplo6/producer_positions.py

```



Animación 7.1_Kafka_Spark: Ejemplo 6

2.7 Persistencia de datos

Si queremos realizar una persistencia de estos datos, sólo tendríamos que almacenarlos **usando algunas de las ya explicadas a lo largo del módulo**, o cualquier otra que sea conveniente para el escenario en concreto (HDFS, Sistema de ficheros, S3, MongoDB,...).

2.8 Conclusiones

A partir de esta configuración, podemos añadir ETL para obtener más conocimiento sobre los jugadores, como velocidades, sprints, mapas de calor, etc, que podría ayudar al stack técnico en diferentes tomas de decisiones (sistema de juego, rendimiento de jugadores, cansancio acumulado,...).

Podríamos añadir información adicional si pudiéramos capturar los datos del balón, que combinados con la de los jugadores, podríamos obtener otros datos valiosos de aplicación de business intelligence para ofrecérselos al stack técnico, tales como: estadísticas de pases, disparos, zonas de mayor juego,...