

Sistemi za obradu i analizu velike količine podataka - Drugi projekat

Tehnologije

Java 8

Hadoop 2.7

Spark 2.4.3

Kafka 2.4.0 (Scala 2.11)

Cassandra DB

Docker

Streaming Producer

```
KafkaProducer<String, String> producer = configureProducer(kafkaUrl);

Configuration conf = new Configuration();
conf.set(CommonConfigurationKeys.FS_DEFAULT_NAME_KEY, hdfsUrl);
conf.set("fs.hdfs.impl", DistributedFileSystem.class.getName());
conf.set("fs.file.impl", LocalFileSystem.class.getName());
FSDataInputStream inputStream = null;
FileSystem fs = null;
try {
    fs = FileSystem.get(URI.create(hdfsUrl), conf);
    Path inFile = new Path(csvFilePath);
    inputStream = fs.open(inFile);
    if (!fs.exists(inFile)) {
        System.out.println("Input file not found");
        throw new IOException("Input file not found");
    }

    String line = inputStream.readLine();
    while (line != null) {
        EventData tmp = EventData.CreateEventData(line);

        if (tmp != null) {
            ProducerRecord<String, String> rec = new ProducerRecord<String, String>(TaxiTopic, line);
            producer.send(rec);
            System.out.println("[KAFKA TAXI DATA SENT]: " + tmp.getVendorId());
            Thread.sleep(dataSendingSleep * 1000);
            System.out.println("Sleeping " + dataSendingSleep + "sec");
        }

        line = inputStream.readLine();
    }
} finally {
    if (inputStream != null) {
        inputStream.close();
    }
    if (fs != null) {
        fs.close();
    }
}

producer.close();
System.out.println("All done.");

System.exit(1);
}
```

Streaming consumer

```
SparkConf conf = new SparkConf().setAppName("BigData-2").setMaster(sparkMasterUrl);
JavaStreamingContext streamingContext = new JavaStreamingContext(conf, new Duration(dataReceivingSleep * 1000));

Map<String, Object> kafkaParams = getKafkaParams(kafkaUrl);
Collection<String> topics = Collections.singletonList(TaxiTopic);

// Creates Kafka Direct Stream.

JavaInputDStream<ConsumerRecord<Object, String>> stream = KafkaUtils.createDirectStream(streamingContext,
    LocationStrategies.PreferConsistent(), ConsumerStrategies.Subscribe(topics, kafkaParams));

JavaDStream<String> receivedData = stream.map(ConsumerRecord::value);
JavaDStream<EventData> eventData = receivedData.map(EventData::CreateEventData);

// Prints received and serialized event data.

eventData.print();

Double minLongitude = -75.971770000000003;
Double maxLongitude = -71.995770000000003;
Double minLatitude = 38.706828000000003;
Double maxLatitude = 42.756998000000003;

// Filters trips with one passenger at a given location.

JavaDStream<EventData> filteredData = eventData.filter(ed -> ed != null &&
    ed.getPickupLongitude() >= minLongitude && ed.getPickupLongitude() <= maxLongitude &&
    ed.getDropoffLongitude() >= minLongitude && ed.getDropoffLongitude() <= maxLongitude &&
    ed.getPickupLatitude() >= minLatitude && ed.getPickupLatitude() <= maxLatitude &&
    ed.getDropoffLatitude() >= minLatitude && ed.getDropoffLatitude() <= maxLatitude &&
    ed.getPassengerCount().equals("1"));
```

Cassandra connector

```
public class CassandraConnector {  
  
    private CqlSession session;  
    private static CassandraConnector instance;  
  
    public static CassandraConnector getInstance() {  
        if (instance == null) {  
            instance = new CassandraConnector();  
        }  
        return instance;  
    }  
  
    public void connect(String addr, Integer port) {  
        CqlSessionBuilder builder = CqlSession.builder();  
        builder.withAuthCredentials("cassandra", "cassandra")  
            .withLocalDatacenter("datacenter1")  
            .addContactPoint(new InetSocketAddress(addr, port));  
  
        session = builder.build();  
    }  
  
    public CqlSession getSession() {  
        return session;  
    }  
  
    public void close() {  
        session.close();  
    }  
  
    private CassandraConnector() {  
    }  
}
```

Prepare Cassandra Keyspace

```
CassandraConnector conn = CassandraConnector.getInstance();
conn.connect(addr, port);
CqlSession session = conn.getSession();

String createKeyspaceCQL = String.format("CREATE KEYSPACE IF NOT EXISTS %s WITH replication "
    + " = {'class':'SimpleStrategy', 'replication_factor':1};", Keyspace);

//String dropTaxiRecordTable = String.format("DROP TABLE %s.%s;", Keyspace, TaxiRecordsTable);

String createTaxiRecordTable = String.format(
    "CREATE TABLE IF NOT EXISTS %s.%s ("
        + " start_time timestamp PRIMARY KEY,"
        + " end_time timestamp,"
        + " fare_amount float,"
        + " pickup_latitude double,"
        + " pickup_longitude double,"
        + " mta_tax float,"
        + " rate_code text,"
        + " surcharge float );",
    Keyspace, TaxiRecordsTable);

//String dropStatisticTable = String.format("DROP TABLE %s.%s;", Keyspace, StatisticsTable);

String createStatisticTable = String.format(
    "CREATE TABLE IF NOT EXISTS %s.%s ("
        + " min_pickup_time timestamp PRIMARY KEY,"
        + " max_pickup_time timestamp,"
        + " min_tip_amount float,"
        + " max_tip_amount float,"
        + " average_tip_amount float,"
        + " trip_count bigint,"
        + " top_payment_type_count int,"
        + " top_payment_type_name text,"
        + " top_rate_code_count int,"
        + " top_rate_code_name text );",
    Keyspace, StatisticsTable);

System.out.println("Preparing Cassandra Keyspace.");

session.execute(createKeyspaceCQL);
System.out.println("Keyspace created");

//session.execute(dropTaxiRecordTable);
//System.out.println(String.format("Table %s dropped.", TaxiRecordsTable));

session.execute(createTaxiRecordTable);
System.out.println(String.format("Table %s created.", StatisticsTable));

//session.execute(dropStatisticTable);
//System.out.println(String.format("Table %s dropped.", StatisticsTable));

session.execute(createStatisticTable);
System.out.println(String.format("Table %s created.", TaxiRecordsTable));
```

Obrada Streaming podataka

```
filteredData.foreachRDD(rdd -> {  
    Long count = rdd.count();  
    if (count <= 0) {  
        System.out.println("INFO: There is no RDD over the filtered data, skipping.");  
        return;  
    }  
  
    // Finds min, max and avg tip amount in current rdd batch.  
  
    EventData minTipAmount = rdd.min(new TipAmountComparator());  
    EventData maxTipAmount = rdd.max(new TipAmountComparator());  
  
    EventData minPickupDate = rdd.min(new DateComparator());  
    EventData maxPickupDate = rdd.max(new DateComparator());  
  
    JavaRDD<Float> tipAmounts = rdd.map((tr) -> tr.getTipAmount());  
  
    Float tipAmountSum = tipAmounts.reduce((tipAmount, accumulator) -> {  
        return tipAmount + accumulator;  
    });  
  
    Float averageTipAmount = tipAmountSum / count;  
  
    // Finds most used payment types in current rdd batch.  
  
    JavaPairRDD<String, Integer> topPaymentTypes = rdd.mapToPair(r -> new Tuple2<>(r.getPaymentType(),1)).reduceByKey(Integer::sum);  
    Integer mostPaymentTypeCount = topPaymentTypes.map(t -> t._2).max(Comparator.naturalOrder());  
  
    List<String> topPaymentTypesNames = topPaymentTypes.filter(t -> t._2.equals(mostPaymentTypeCount)).map(t -> t._1).collect();  
    String mostPaymentTypeNames = String.join(",", topPaymentTypesNames);  
  
    // Finds most used rate codes in current rdd batch.  
  
    JavaPairRDD<String, Integer> topRateCodes = rdd.mapToPair(r -> new Tuple2<>(r.getRateCode(),1)).reduceByKey(Integer::sum);  
    Integer mostRateCodesCount = topRateCodes.map(t -> t._2).max(Comparator.naturalOrder());  
  
    List<String> topRateCodeNames = topRateCodes.filter(t -> t._2.equals(mostRateCodesCount)).map(t -> t._1).collect();  
    String mostRateCodeNames = String.join(",", topRateCodeNames);  
  
    // Creates and stores statistic data in Cassandra Db.  
  
    StatisticRecord statisticData = new StatisticRecord(minPickupDate.getPickupDateTime(), maxPickupDate.getPickupDateTime(),  
        minTipAmount.getTipAmount(), maxTipAmount.getTipAmount(), averageTipAmount, count,  
        mostPaymentTypeCount, mostPaymentTypeNames, mostRateCodesCount, mostRateCodeNames);  
  
    //System.out.println("Statistic Data");  
    //System.out.println(statisticData.toString());  
    saveStatisticRecord(statisticData, cassandraUrl, cassandraPort);  
    System.out.println("INFO: Statistics are stored in database.");  
}
```

Smeštanje taxi rekorda u Cassandra DB

```
// Stores some record attributes in Cassandra Db for testing purpose.  
rdd.foreach(c -> {  
    TaxiRecord taxiRecord = new TaxiRecord(c.getPickupDateTime(), c.getDropoffDateTime(),  
        c.getFareAmount(), c.getPickupLatitude(),  
        c.getPickupLongitude(), c.getMtaTax(),  
        c.getRateCode(), c.getSurcharge());  
    saveTaxiRecord(taxiRecord, cassandraUrl, cassandraPort);  
    System.out.println("INFO: Taxi records are stored in database.");  
});  
  
// Reading data from Cassandra Db.  
  
// readRecords(StatisticsTable, 10, cassandraUrl, cassandraPort);  
// readRecords(TaxiRecordsTable, 10, cassandraUrl, cassandraPort);  
});  
  
streamingContext.start();  
streamingContext.awaitTermination();
```


Smeštanje podataka u Cassandra Db

```
public static void saveStatisticRecord(StatisticRecord statisticRecord, String addr, Integer port) {  
  
    CassandraConnector conn = CassandraConnector.getInstance();  
    conn.connect(addr, port);  
    CqlSession session = conn.getSession();  
  
    RegularInsert insertInto = QueryBuilder.insertInto(Keyspace, StatisticsTable)  
        .value("min_pickup_time", QueryBuilder.bindMarker())  
        .value("max_pickup_time", QueryBuilder.bindMarker())  
        .value("min_tip_amount", QueryBuilder.bindMarker())  
        .value("max_tip_amount", QueryBuilder.bindMarker())  
        .value("average_tip_amount", QueryBuilder.bindMarker())  
        .value("trip_count", QueryBuilder.bindMarker())  
        .value("top_payment_type_count", QueryBuilder.bindMarker())  
        .value("top_payment_type_name", QueryBuilder.bindMarker())  
        .value("top_rate_code_count", QueryBuilder.bindMarker())  
        .value("top_rate_code_name", QueryBuilder.bindMarker());  
  
    SimpleStatement insertStatement = insertInto.build();  
    PreparedStatement preparedStatement = session.prepare(insertStatement);  
  
    BoundStatement boundStatement = preparedStatement.bind()  
        .setInstant(0, statisticRecord.getMinPickupTime().toInstant())  
        .setInstant(1, statisticRecord.getMaxPickupTime().toInstant())  
        .setFloat(2, statisticRecord.getMinTipAmount())  
        .setFloat(3, statisticRecord.getMaxTipAmount())  
        .setFloat(4, statisticRecord.getAverageTipAmount())  
        .setLong(5, statisticRecord.getTripCount())  
        .setInt(6, statisticRecord.getTopPaymentTypeCount())  
        .setString(7, statisticRecord.getTopPaymentTypeName())  
        .setInt(8, statisticRecord.getTopRateCodeCount())  
        .setString(9, statisticRecord.getTopRateCodeName());  
  
    session.execute(boundStatement);  
    System.out.println(String.format("Data stored in %s table.", TaxiRecordsTable));  
  
    session.close();  
    conn.close();  
}
```

```
public static void saveTaxiRecord(TaxiRecord taxiRecord, String addr, Integer port) {  
  
    CassandraConnector conn = CassandraConnector.getInstance();  
    conn.connect(addr, port);  
    CqlSession session = conn.getSession();  
  
    RegularInsert insertInto = QueryBuilder  
        .insertInto(Keyspace, TaxiRecordsTable)  
        .value("start_time", QueryBuilder.bindMarker())  
        .value("end_time", QueryBuilder.bindMarker())  
        .value("fare_amount", QueryBuilder.bindMarker())  
        .value("pickup_latitude", QueryBuilder.bindMarker())  
        .value("pickup_longitude", QueryBuilder.bindMarker())  
        .value("mta_tax", QueryBuilder.bindMarker())  
        .value("rate_code", QueryBuilder.bindMarker())  
        .value("surcharge", QueryBuilder.bindMarker());  
  
    SimpleStatement insertStatement = insertInto.build();  
    PreparedStatement preparedStatement = session.prepare(insertStatement);  
  
    BoundStatement boundStatement = preparedStatement.bind()  
        .setInstant(0, taxiRecord.getStartTime().toInstant())  
        .setInstant(1, taxiRecord.getEndTime().toInstant())  
        .setFloat(2, taxiRecord.getFareAmount())  
        .setDouble(3, taxiRecord.getPickupLatitude())  
        .setDouble(4, taxiRecord.getPickupLongitude())  
        .setFloat(5, taxiRecord.getMtaTax())  
        .setString(6, taxiRecord.getRateCode())  
        .setFloat(7, taxiRecord.getSurcharge());  
  
    session.execute(boundStatement);  
    conn.close();  
}
```

Statistic Table (TablePlus DB Management alat)

<

Taxi Record Table (TablePlus DB Management alat)

Cassandra 3.11.10 : Big-Data : test : taxi_records

statistics

Search for item...

Items Queries History

statistics

taxi_records

Unset SQL Export

Show: Ctrl+F Hide: ESC Insert: Ctrl+I Remove: Ctrl+Shift+I Up: Ctrl+I Down: Ctrl+I Apply: Enter Apply all: Ctrl+Enter

start_time	end_time	fare_amount	mta_tax	pickup_latitude	pickup_longitude	rate_code	surcharge
2014-01-09 19:06:34	2014-01-09 19:12:45	7.5	0.5	-73.938787	-73.938787	1	1
2014-01-09 19:25:25	2014-01-09 19:36:00	9.5	0.5	-73.968523	-73.968523	1	1
2014-01-09 20:25:48	2014-01-09 20:36:47	9	0.5	-73.990014	-73.990014	1	0.5
2014-01-09 19:49:48	2014-01-09 19:58:41	7.5	0.5	-73.991955	-73.991955	1	1
2014-01-09 23:39:39	2014-01-09 23:45:20	7	0.5	-73.95275	-73.95275	1	0.5
2014-01-09 22:01:46	2014-01-09 22:14:01	9.5	0.5	-73.987331	-73.987331	1	0.5
2014-01-09 20:42:45	2014-01-09 20:46:52	4.5	0.5	-73.993591	-73.993591	1	0.5
2014-01-09 21:54:18	2014-01-09 22:04:39	10.5	0.5	-73.981957	-73.981957	1	0.5
2014-01-09 17:42:05	2014-01-09 17:47:04	6	0.5	-73.987374	-73.987374	1	1
2014-01-09 18:43:51	2014-01-09 19:16:34	30	0.5	-73.950583	-73.950583	1	1
2014-01-09 23:27:52	2014-01-09 23:46:30	16	0.5	-73.994892	-73.994892	1	0.5
2014-01-09 20:59:19	2014-01-09 21:10:51	10.5	0.5	-73.979001	-73.979001	1	0.5
2014-01-09 19:38:54	2014-01-09 19:45:10	6	0.5	-74.002372	-74.002372	1	1
2014-01-09 20:24:42	2014-01-09 20:43:12	16.5	0.5	-74.000847	-74.000847	1	0.5
2014-01-09 21:55:58	2014-01-09 22:10:48	13.5	0.5	-73.982095	-73.982095	1	0.5
2014-01-09 23:59:53	2014-01-10 00:06:38	8	0.5	-73.98972	-73.98972	1	0.5
2014-01-09 20:59:31	2014-01-09 21:13:56	10.5	0.5	-73.987459	-73.987459	1	0.5
2014-01-09 17:13:03	2014-01-09 18:12:57	50	0.5	-73.994028	-73.994028	1	1
2014-01-09 20:18:12	2014-01-09 20:36:14	15	0.5	-74.003451	-74.003451	1	0.5
2014-01-09 20:09:16	2014-01-09 21:03:10	52	0.5	-73.988093	-73.988093	2	0
2014-01-09 20:13:40	2014-01-09 20:24:06	10.5	0.5	-73.980528	-73.980528	1	0.5
2014-01-09 20:19:01	2014-01-09 20:22:14	4	0.5	-74.00496	-74.00496	1	0.5
2014-01-09 17:49:07	2014-01-09 17:50:46	4	0.5	-74.002374	-74.002374	1	1
2014-01-09 22:35:05	2014-01-09 22:48:03	12	0.5	-73.989117	-73.989117	1	0.5
2014-01-09 19:38:31	2014-01-09 19:53:37	12	0.5	-73.993447	-73.993447	1	1
2014-01-09 21:55:45	2014-01-09 22:03:56	8.5	0.5	-73.977888	-73.977888	1	0.5
2014-01-09 16:45:21	2014-01-09 16:49:01	4.5	0.5	-73.977261	-73.977261	1	1
2014-01-09 17:33:41	2014-01-09 17:39:04	5.5	0.5	-73.982441	-73.982441	1	1
2014-01-09 21:53:25	2014-01-09 22:01:52	8.5	0.5	-73.975407	-73.975407	1	0.5
2014-01-09 19:31:22	2014-01-09 19:56:28	23.5	0.5	-73.970811	-73.970811	1	1
2014-01-09 19:54:09	2014-01-09 20:25:07	52	0.5	-73.783365	-73.783365	2	0
2014-01-09 20:17:23	2014-01-09 20:38:44	16.5	0.5	-73.968604	-73.968604	1	0.5
2014-01-09 21:07:16	2014-01-09 21:10:51	4.5	0.5	-73.984586	-73.984586	1	0.5
2014-01-09 20:13:11	2014-01-09 20:18:53	5.5	0.5	-73.992874	-73.992874	1	0.5
2014-01-09 20:50:24	2014-01-09 21:00:13	8.5	0.5	-73.96147	-73.96147	1	0.5

1-144 of 0 rows

Columns Filters

Pomoćne funkcije

```
public static void readRecords(String tableName, int limit, String addr, Integer port) {  
    CassandraConnector conn = CassandraConnector.getInstance();  
    conn.connect(addr, port);  
    CqlSession session = conn.getSession();  
  
    String query = String.format("SELECT * FROM %s.%s limit %d" ,Keyspace, tableName, limit);  
    ResultSet result = session.execute(query);  
  
    System.out.println(String.format("First %d records from table %s", limit, tableName));  
    System.out.println(result.all());  
  
    session.close();  
    conn.close();  
}
```

```
public static Map<String, Object> getKafkaParams(String brokers) {  
    Map<String, Object> kafkaParams = new HashMap<String, Object>();  
    kafkaParams.put(ConsumerConfig.BOOTSTRAP_SERVERS_CONFIG, brokers);  
    kafkaParams.put(ConsumerConfig.CLIENT_ID_CONFIG, "streaming-consumer");  
    kafkaParams.put(ConsumerConfig.GROUP_ID_CONFIG, "streaming-consumer");  
    kafkaParams.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG, "earliest");  
    kafkaParams.put(ConsumerConfig.KEY_DESERIALIZER_CLASS_CONFIG, StringDeserializer.class);  
    kafkaParams.put(ConsumerConfig.VALUE_DESERIALIZER_CLASS_CONFIG, StringDeserializer.class);  
    kafkaParams.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG, "latest");  
    kafkaParams.put(ConsumerConfig.ENABLE_AUTO_COMMIT_CONFIG, false);  
  
    return kafkaParams;  
}
```