# Sistemi za obradu i analizu velike količine podataka - Treći projekat

Nikola Božinović 1030

Tehnologije

- Java 8
- Hadoop 2.7
- Spark 2.4.3
- Kafka 2.4.0 (Scala 2.11)
- Docker

## ML Batch

```java
SparkSession spark = SparkSession.builder().appName("BigData-3-ML-Saving").master(sparkMasterUrl).getOrCreate();

Dataset<Row> dataSet = spark.read().option("header", "true").csv(csvFile);

Dataset<Row> filteredData = dataSet.filter((row) -> {
    return !row.anyNull();
});

UDF1<String, Integer> udfGetDayOfWeak = Main::getDayOfWeekFromDate;
UserDefinedFunction getDayOfWeak = functions.udf(udfGetDayOfWeak, DataTypes.IntegerType);

UDF1<String, Integer> udfGetHourOfDay = Main::getHourOfDayFromDate;
UserDefinedFunction getHourOfDay = functions.udf(udfGetHourOfDay, DataTypes.IntegerType);

Dataset<Row> selectedData = filteredData.select(
        dataSet.col("trip_distance").cast(FloatType).as("TripDistance"),
        dataSet.col("pickup_longitude").cast(DoubleType).as("PickupLongitude"),
        dataSet.col("pickup_latitude").cast(DoubleType).as("PickupLatitude"),
        dataSet.col("dropoff_longitude").cast(DoubleType).as("DropoffLongitude"),
        dataSet.col("dropoff_latitude").cast(DoubleType).as("DropoffLatitude"),
        dataSet.col("rate_code").cast(IntegerType).as("RateCode"),
        getDayOfWeak.apply(dataSet.col("pickup_datetime")).as("DayOfWeak"),
        getHourOfDay.apply(dataSet.col("pickup_datetime")).as("HourOfDay"),
        dataSet.col("fare_amount").cast(FloatType).as("FareAmount")
        );
```

## ML Batch

```java
// nyc coordinates:
// latitude = 40.730610
// longitude = -73.935242
double nycMinLatitude = 38.730610; // -2
double nycMaxLatitude = 42.730610; // +2
double nycMinLongitude = -75.935242; // -2
double nycMaxLongitude = -71.935242; // +2

Dataset<Row> filteredAndSelectedData = selectedData.filter(selectedData.col("TripDistance").notEqual(0.0)
        .and(selectedData.col("PickupLongitude").gt(nycMinLongitude))
        .and(selectedData.col("PickupLongitude").lt(nycMaxLongitude))
        .and(selectedData.col("PickupLatitude").gt(nycMinLatitude))
        .and(selectedData.col("PickupLatitude").lt(nycMaxLatitude))
        .and(selectedData.col("DropoffLongitude").gt(nycMinLongitude))
        .and(selectedData.col("DropoffLongitude").lt(nycMaxLongitude))
        .and(selectedData.col("DropoffLatitude").gt(nycMinLatitude))
        .and(selectedData.col("DropoffLatitude").lt(nycMaxLatitude)));


VectorAssembler vectorAssembler = new VectorAssembler()
        .setInputCols(new String[]{"TripDistance", "PickupLongitude", "PickupLatitude", "DropoffLongitude",
                "DropoffLatitude","DayOfWeak", "HourOfDay","RateCode"})
        .setOutputCol("Features");

Dataset<Row> transformedData = vectorAssembler.transform(filteredAndSelectedData);
```
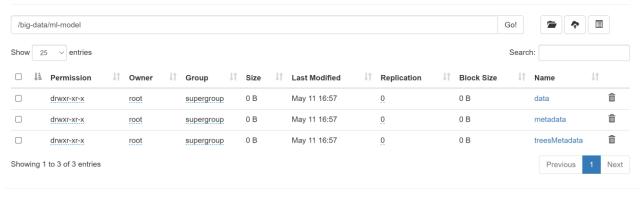
ML Batch

```java
Dataset<Row>[] splits = transformedData.randomSplit(new double[]{0.7, 0.3});
Dataset<Row> trainingData = splits[0];
Dataset<Row> testData = splits[1];

RandomForestRegressor rf = new RandomForestRegressor()
        .setLabelCol("FareAmount")
        .setFeaturesCol("Features");

RandomForestRegressionModel model = rf.fit((trainingData));

model.write().overwrite().save(hdfsUrl + "/big-data/ml-model");

Dataset<Row> predictions = model.transform(testData);
predictions.show(100);

RegressionEvaluator evaluator = new RegressionEvaluator()
        .setLabelCol("FareAmount")
        .setPredictionCol("prediction")
        .setMetricName("rmse");

double rmse = evaluator.evaluate(predictions);
System.out.println("Root Mean Squared Error (RMSE) on test data = " + rmse);

spark.stop();
spark.close();
```

# Browse Directory

/big-data/ml-model                                Go!

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | drwxr-xr-x | root | supergroup | 0 B | May 11 16:57 | 0 | 0 B | data | 🗑 |
| ☐ | drwxr-xr-x | root | supergroup | 0 B | May 11 16:57 | 0 | 0 B | metadata | 🗑 |
| ☐ | drwxr-xr-x | root | supergroup | 0 B | May 11 16:57 | 0 | 0 B | treesMetadata | 🗑 |

Showing 1 to 3 of 3 entries

Previous   1   Next

Hadoop, 2019.

**ML Batch**

/big-data/ml-model/data                                Go!

Show 25 entries                                                    Search:

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | root | supergroup | 0 B | May 11 16:57 | 3 | 128 MB | _SUCCESS | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.78 KB | May 11 16:57 | 3 | 128 MB | part-00000-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.81 KB | May 11 16:57 | 3 | 128 MB | part-00001-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.85 KB | May 11 16:57 | 3 | 128 MB | part-00002-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 9.54 KB | May 11 16:57 | 3 | 128 MB | part-00003-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.77 KB | May 11 16:57 | 3 | 128 MB | part-00004-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.83 KB | May 11 16:57 | 3 | 128 MB | part-00005-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.81 KB | May 11 16:57 | 3 | 128 MB | part-00006-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 9.85 KB | May 11 16:57 | 3 | 128 MB | part-00007-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.84 KB | May 11 16:57 | 3 | 128 MB | part-00008-3f1c9c4a-693b-43a6-a83a-fd88062f01f4-c000.snappy.parquet | 🗑 |
| ☐ | -rw-r--r-- | root | supergroup | 6.83 KB | May 11 16:57 | 3 | 128 MB | part-00009-3f1c9c4a-693b-43a6-a83a-fd88062f01f4- | 🗑 |

## ML Streaming

```java
SparkSession spark = SparkSession.builder().appName("BigData-4-ML-Streaming").master(sparkMasterUrl).getOrCreate();
JavaSparkContext javaSparkContext = JavaSparkContext.fromSparkContext(spark.sparkContext());
JavaStreamingContext streamingContext = new JavaStreamingContext(javaSparkContext, new Duration(dataReceivingSleep * 1000)

RandomForestRegressionModel model = RandomForestRegressionModel.load(hdfsUrl + "/big-data/ml-model/");

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

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);

JavaDStream<EventData> filteredData = eventData.filter(ed -> ed != null &&
        !ed.getTripDistance().equals("0.0") &&
        !ed.getPickupLongitute().equals("0.0") &&
        !ed.getPickupLatitude().equals("0.0") &&
        !ed.getDropoffLongitude().equals("0.0") &&
        !ed.getDropoffLatitude().equals("0.0") &&
        !ed.getRateCode().equals(null) &&
        !ed.getPickupDateTime().equals(null)
);

JavaDStream<Row> rows = filteredData.map(row -> RowFactory.create(
        Float.parseFloat(row.getTripDistance()),
        Double.parseDouble(row.getPickupLongitute()),
        Double.parseDouble(row.getPickupLatitude()),
        Double.parseDouble(row.getDropoffLongitude()),
        Double.parseDouble(row.getDropoffLatitude()),
        Integer.parseInt(row.getRateCode()),
        getDayOfWeekFromDate(row.getPickupDateTime()),
        getHourOfDayFromDate(row.getPickupDateTime())));
```

# ML Streaming

```java
rows.foreachRDD(d -> {
    StructType rowSchema = DataTypes.createStructType(
        new StructField[]{
            createStructField("TripDistance", DataTypes.FloatType, false),
            createStructField("PickupLongitude", DataTypes.DoubleType, false),
            createStructField("PickupLatitude", DataTypes.DoubleType, false),
            createStructField("DropoffLongitude", DataTypes.DoubleType, false),
            createStructField("DropoffLatitude", DataTypes.DoubleType, false),
            createStructField("RateCode", DataTypes.IntegerType, false),
            createStructField("DayOfWeak", DataTypes.IntegerType, false),
            createStructField("HourOfDay", DataTypes.IntegerType, false),
    });

    Dataset<Row> data = spark.createDataFrame(d, rowSchema);

    VectorAssembler vectorAssembler = new VectorAssembler()
        .setInputCols(new String[]{"TripDistance", "PickupLongitude", "PickupLatitude", "DropoffLongitude",
            "DropoffLatitude", "RateCode", "DayOfWeak", "HourOfDay"})
        .setOutputCol("Features");

    Dataset<Row> transformed = vectorAssembler.transform(data);
    Dataset<Row> predictions = model.transform(transformed);
    predictions.show(100);
});
```

## Pomoćne funkcije

```java
public static int getDayOfWeekFromDate(String stringDate) {
    // example format: 2014-01-09 20:45:25
    SimpleDateFormat formatter = new SimpleDateFormat("yyyy-MM-dd HH:mm:ss");
    Date date = new Date();

    try {
        date = formatter.parse(stringDate);
    } catch (ParseException e) {
        e.printStackTrace();
    }

    Calendar calendar = Calendar.getInstance();
    calendar.setTime(date);
    return calendar.get(Calendar.DAY_OF_WEEK); // the day of the week in numerical format
}

public static int getHourOfDayFromDate(String stringDate) {
    // example format: 2014-01-09 20:45:25
    SimpleDateFormat formatter = new SimpleDateFormat("yyyy-MM-dd HH:mm:ss");
    Date date = new Date();

    try {
        date = formatter.parse(stringDate);
    } catch (ParseException e) {
        e.printStackTrace();
    }

    Calendar calendar = Calendar.getInstance();
    calendar.setTime(date);
    return calendar.get(Calendar.HOUR_OF_DAY); // the hour of the day in numerical format
}
```