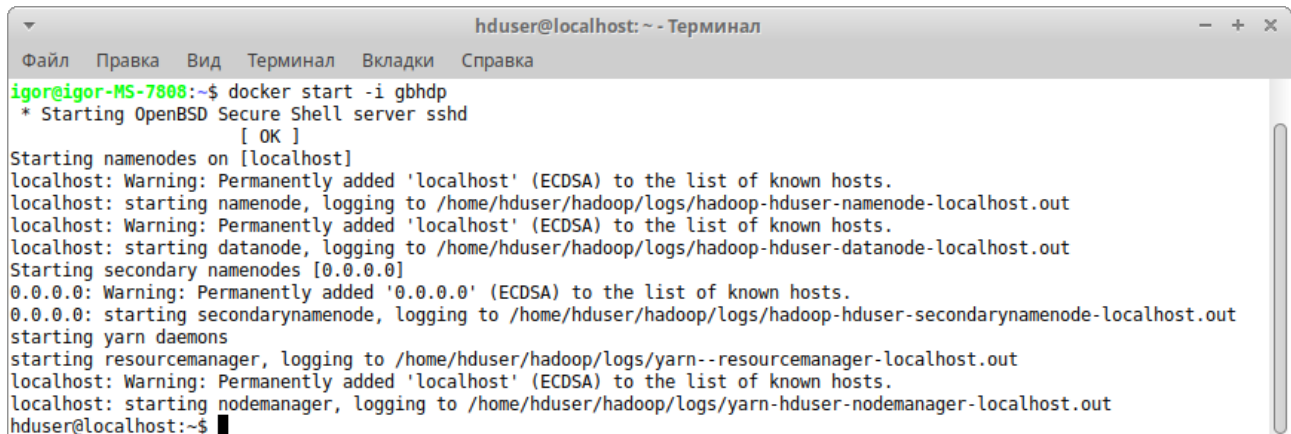1. Проведите анализ тональности датасета [IMDB](): обучите модель на Train.csv, после чего проверьте её на Valid.csv

2. Посчитайте получившуюся точность (accuracy — количество правильных предсказаний от общего количества ответов) модели. Ваш код должен использовать ML Pipelines и использовать трансформеры для подготовки данных к обучению
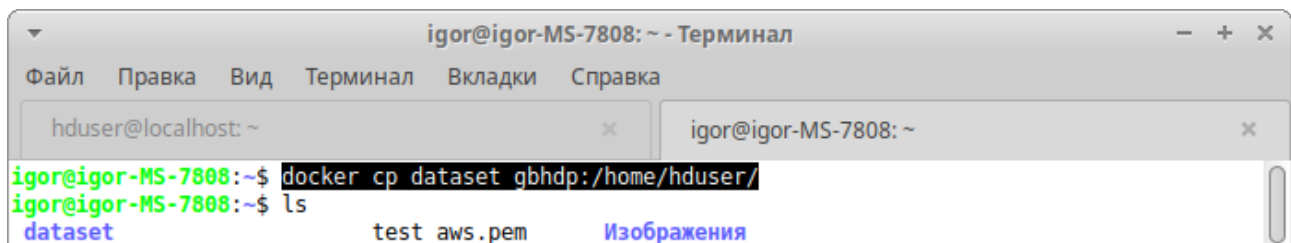
Стартуем остановленный контейнер с Hadoop и Spark:

$ docker start -i gbhdp



Копируем датасет в контейнер: docker cp dataset gbhdp:/home/hduser/



Копируем директорию в hdfs: hdfs dfs -put dataset /user/hduser/

Проверим, что файлы в hdfs: hdfs dfs -ls /user/hduser/dataset



Копируем датасет в контейнер:

docker cp Train.csv gbhdp:/home/hduser/

docker cp Valid.csv gbhdp:/home/hduser/

Файл   Правка   Вид   Терминал   Вкладки   Справка

| hduser@localhost: ~/gbhdp ✕ | igor@igor-MS-7808: ~/dataset ✕ |

```
igor@igor-MS-7808:~/dataset$ ls
Test.csv  Train.csv  Valid.csv
igor@igor-MS-7808:~/dataset$ docker cp Train.csv gbhdp:/home/hduser/
igor@igor-MS-7808:~/dataset$ docker cp Valid.csv gbhdp:/home/hduser/
igor@igor-MS-7808:~/dataset$ ▮
```

создаём директорию imdb: mkdir imdb

Файл   Правка   Вид   Терминал   Вкладки   Справка

| hduser@localhost: ~ ✕ | igor@igor-MS-7808: ~ ✕ |

```
hduser@localhost:~$ mkdir imdb
hduser@localhost:~$ ls
MapReduce-1.0-SNAPSHOT.jar  Train.csv  hadoop  scalasparkrdd_2.11-0.1.jar
SparkDF-1.0-SNAPSHOT.jar    Valid.csv  imdb    spark
SparkRDD-1.0-SNAPSHOT.jar   dataset    ppkm    sparkdfscala_2.11-0.1.jar
hduser@localhost:~$ ▮
```

hdfs dfs -put Train.csv /user/hduser/imdb/

hdfs dfs -put Valid.csv /user/hduser/imdb/

SCALA датасета IMDB: ПРОЕКТ В IDEA

**New Project** ✕

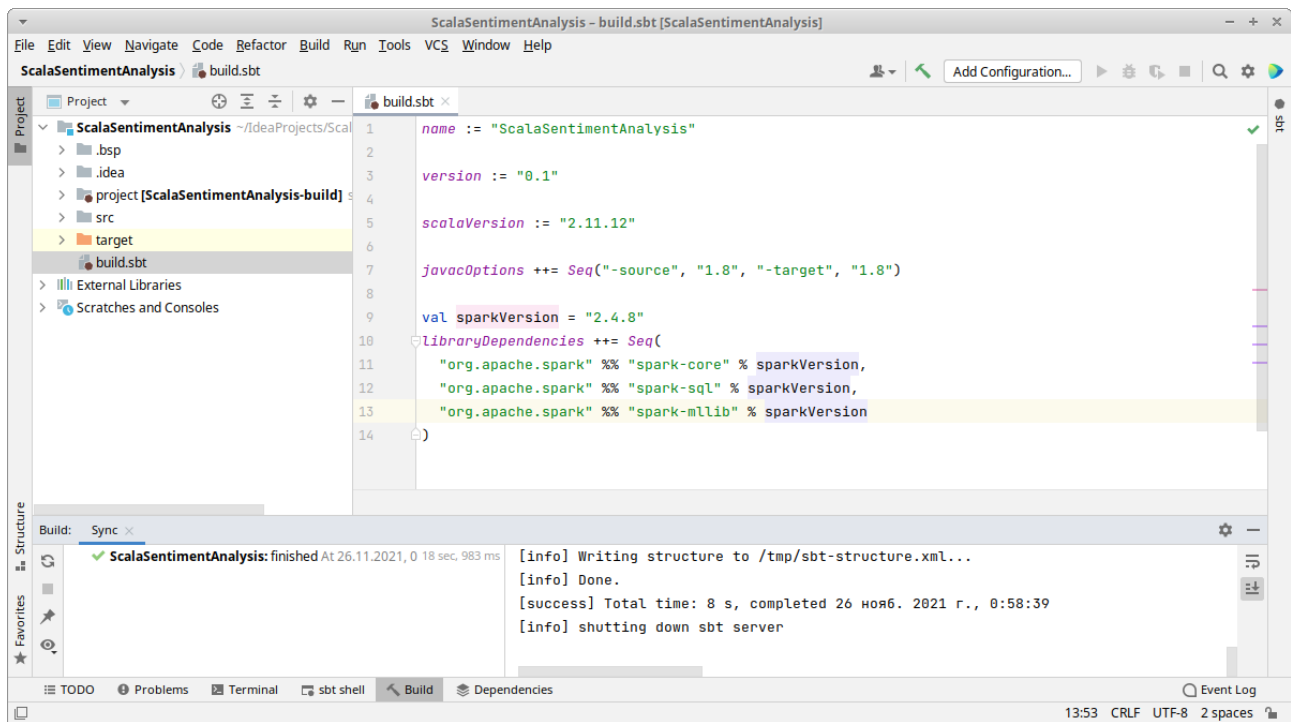| Name: | ScalaSentimentAnalysis |
| Location: | ~/IdeaProjects/ScalaSentimentAnalysis ... |
| JDK: | 🖳 11 version 11.0.11 ▼ |
| sbt: | 1.5.5 ▼  ☐ Sources |
| Scala: | 2.11.12 ▼  ☑ Sources |
| Package prefix: | Such as 'org.example.application' ⑦ |

▸ More Settings

Previous   **Finish**   Cancel   Help

Добавим в build.sbt зависимость на Apache Spark:

Создадим ScalaSentimentAnalysis.scala:

```scala
import org.apache.spark.ml.feature.StopWordsRemover
import org.apache.spark.sql.SparkSession
import org.apache.spark.ml.feature.{HashingTF, StopWordsRemover, Tokenizer}
import org.apache.spark.ml.classification.LogisticRegression
import org.apache.spark.ml.PipelineStage
import org.apache.spark.ml.Pipeline


object SentimentAnalysis {
  def main(args: Array[String]): Unit = {
    val trainCSV = args(0)
    val testCSV = args(1)
    val spark = SparkSession.builder().getOrCreate()

    import spark.sqlContext.implicits._


    // TRAINING
    val trainingDF = spark.read
      .option("header", "true")
      .option("quotes", "\"")
      .option("escape", "\"")
      .option("inferSchema", "true")
      .csv(trainCSV)

    val tokenizer = new Tokenizer()
      .setInputCol("text")
      .setOutputCol("words")

    val remover = new StopWordsRemover()
      .setInputCol("words")
      .setOutputCol("filtered")

    val hashingTF = new HashingTF()
      .setInputCol(remover.getOutputCol)
      .setOutputCol("features")
      .setNumFeatures(1000)

    val lr = new LogisticRegression()
      .setMaxIter(10)
      .setRegParam(0.001)

    val pipeline = new Pipeline()
      .setStages(Array[PipelineStage](tokenizer, remover, hashingTF, lr))

    val model = pipeline.fit(trainingDF)

    // PREDICTION
    val testDF = spark.read
      .option("header", "true")
      .option("quotes", "\"")
      .option("escape", "\"")
      .option("inferSchema", "true")
      .csv(testCSV)

    val rowsCount = testDF.count()

    val predictionDF = model.transform(testDF)

    predictionDF.show( numRows = 20)

    predictionDF
      .groupBy( cols = $"label" - $"prediction" as "result").count()
      .select( cols = $"count" / rowsCount as "accuracy")
      .filter( condition = $"result" === 0)
      .select( cols = $"accuracy")
      .show()

    spark.stop()
  }
}
```

Build:  Sync ×

✔ ScalaSentimentAnalysis: finished At 26.11.2021, 0 18 sec, 983 ms

```
[info] Writing structure to /tmp/sbt-structure.xml...
[info] Done.
[success] Total time: 8 s, completed 26 ноя6. 2021 г., 0:58:39
[info] shutting down sbt server
```

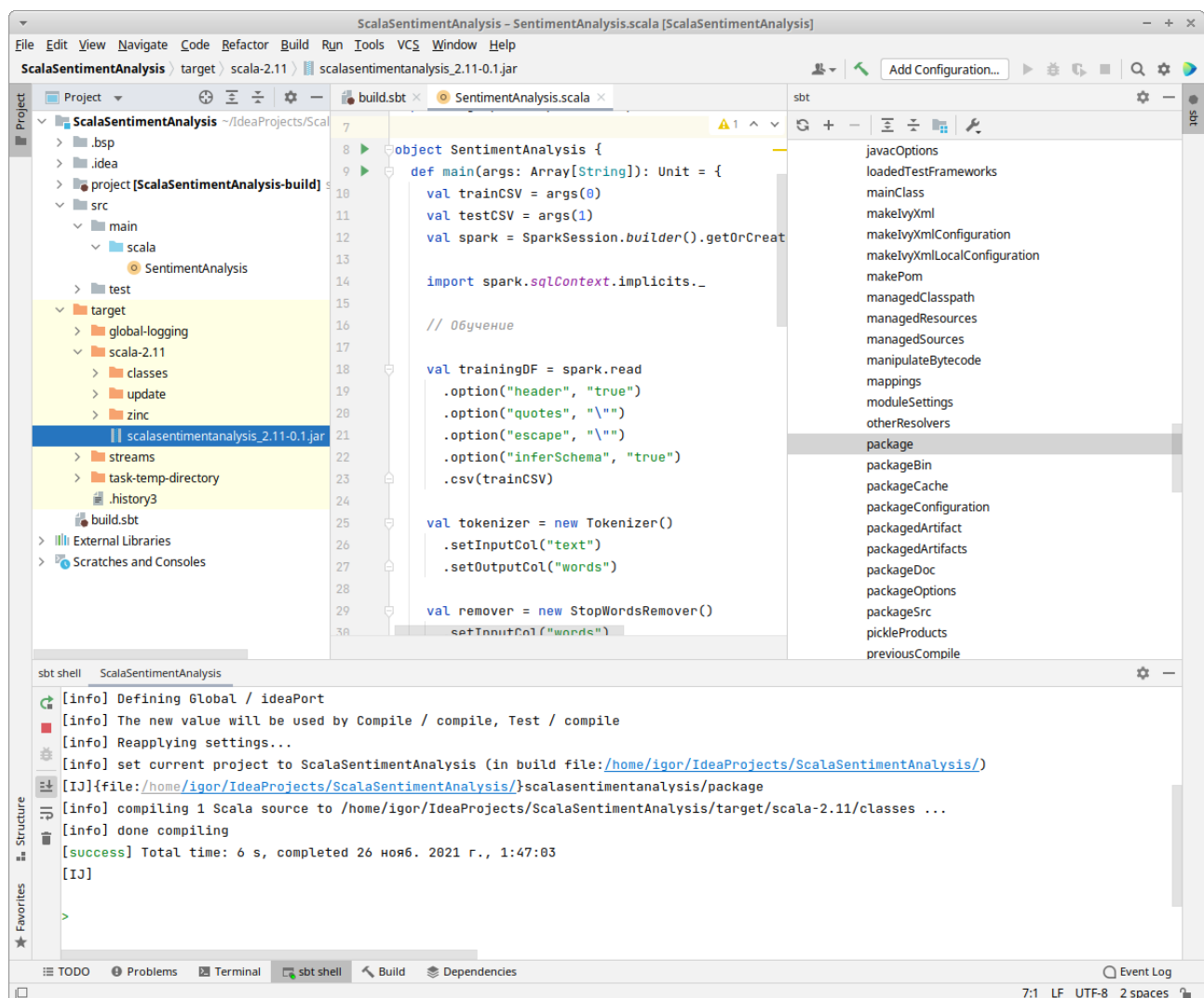TODO   ⊘ Problems   ⌦ Terminal   ⌦ sbt shell   ⚒ Build   ⬥ Dependencies                     ○ Event Log
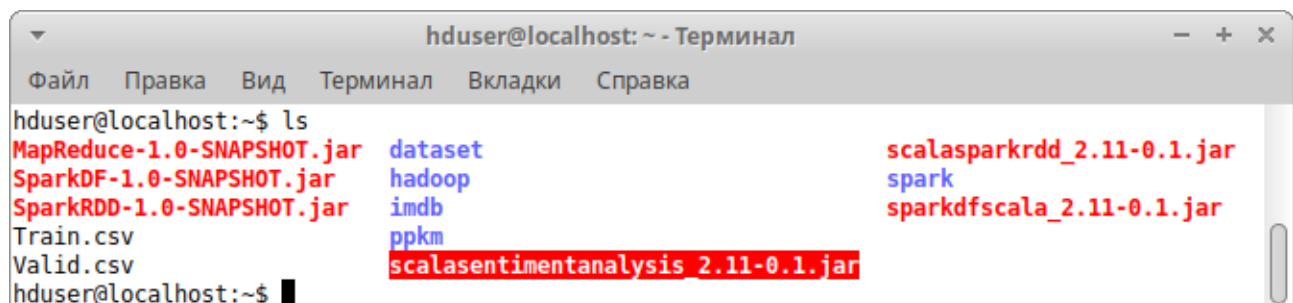
15:1   LF   UTF-8   2 spaces

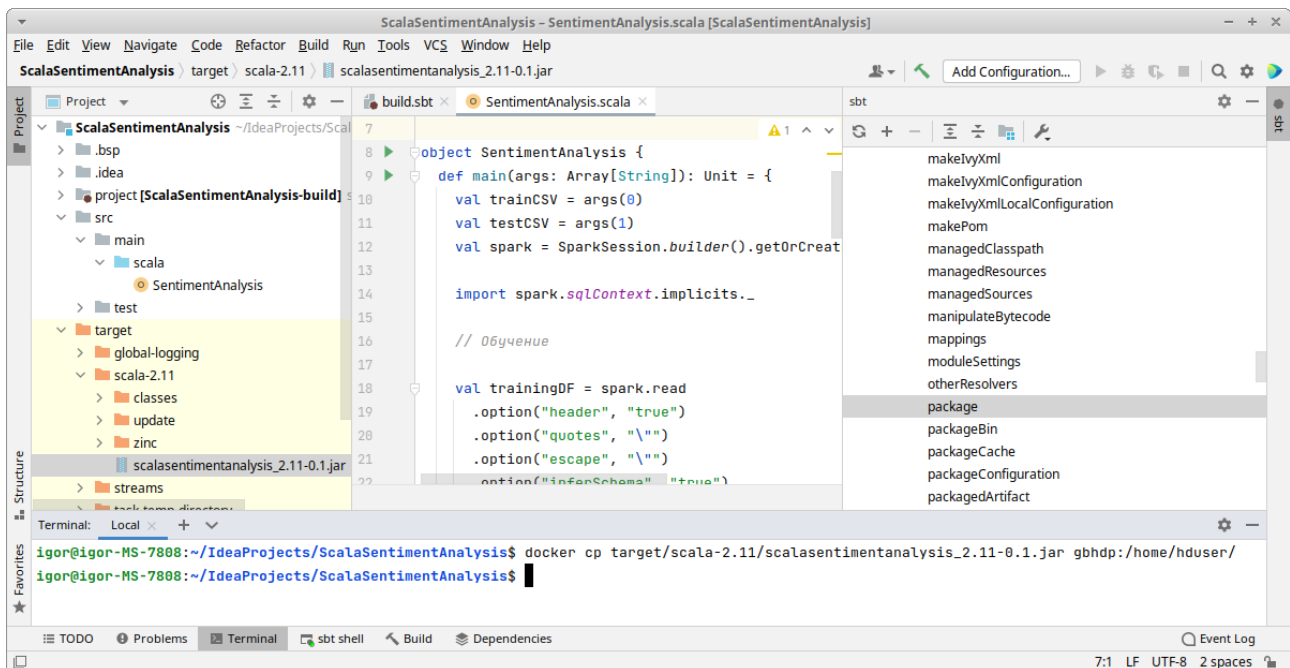# SCALA `SentimentAnalysis`: СБОРКА И ДОСТАВКА: sbt package

Соберем jar: sbt package



Перенесем его на кластер:

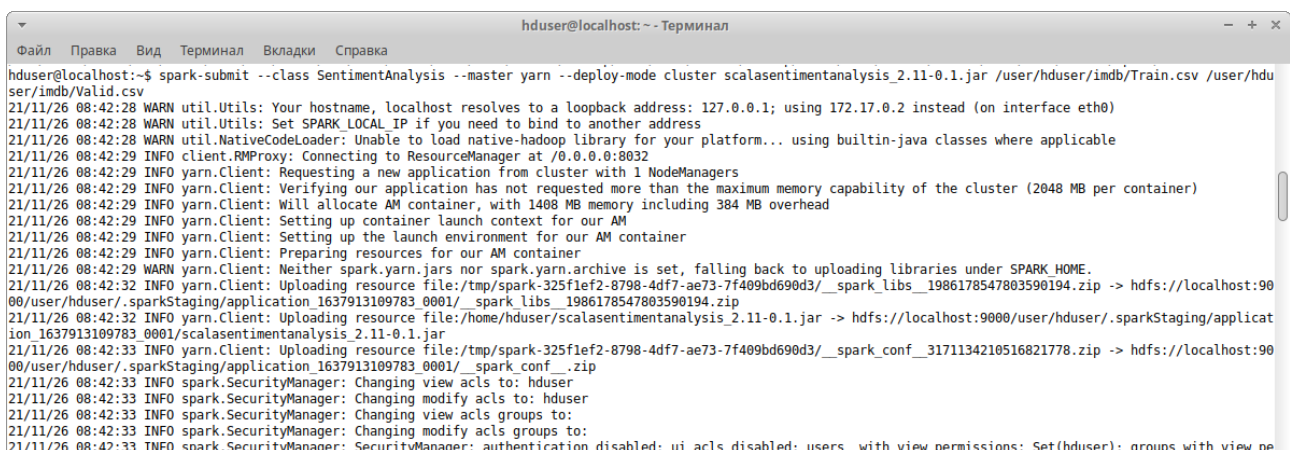docker cp target/scala-2.11/scalasentimentanalysis_2.11-0.1.jar gbhdp:/home/hduser/

Запускаем приложение:

spark-submit --class SentimentAnalysis --master yarn --deploy-mode cluster scalasentimentanalysis_2.11-0.1.jar /user/hduser/imdb/Train.csv /user/hduser/imdb/Valid.csv



tracking URL: http://localhost:8088/proxy/application_1637913109783_0001/

Logs for container_1637913109783_0001_01_000001



accuracy|

+--------+
|  0.7624|