

Ls -la

```
hduser@localhost:~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
hduser@localhost:~$ ls -la
total 76
drwxr-xr-x 1 hduser hduser 4096 Nov 24 19:15 .
drwxr-xr-x 1 root  root  4096 Nov 22 20:40 ..
-rw----- 1 hduser hduser 1535 Nov 24 19:37 .bash_history
-rw-r--r-- 1 hduser hduser 220 Apr  4 2018 .bash_logout
-rw-r--r-- 1 hduser hduser 3889 Nov 24 17:56 .bashrc
drwx----- 2 hduser hduser 4096 Nov 22 20:43 .cache
-rw-r--r-- 1 hduser hduser 807 Apr  4 2018 .profile
drwx----- 2 hduser hduser 4096 Nov 22 20:40 .ssh
-rw-r--r-- 1 hduser hduser    0 Nov 22 20:40 .sudo_as_admin_successful
-rw----- 1 hduser hduser 1124 Nov 24 17:56 .viminfo
-rw-rw-r-- 1 hduser hduser 4441 Nov 22 22:32 MapReduce-1.0-SNAPSHOT.jar
-rw-rw-r-- 1 hduser hduser 2380 Nov 24 18:35 SparkRDD-1.0-SNAPSHOT.jar
drwxr-xr-x 1 hduser hduser 4096 Nov 22 20:43 hadoop
drwxr-xr-x 2 hduser hduser 4096 Jan  1 1970 ppkm
-rw-rw-r-- 1 hduser hduser 6604 Nov 24 19:09 scalasparkrdd_2.11-0.1.jar
drwxr-xr-x 13 hduser hduser 4096 May  8 2021 spark
hduser@localhost:~$
```

tail .bashrc

```
hduser@localhost:~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
hduser@localhost:~$ tail .bashrc
if ! shopt -oq posix; then
  if [ -f /usr/share/bash-completion/bash_completion ]; then
    . /usr/share/bash-completion/bash_completion
  elif [ -f /etc/bash_completion ]; then
    . /etc/bash_completion
  fi
fi
export SPARK_HOME=/home/hduser/spark
export PATH=$PATH:$SPARK_HOME/bin
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
hduser@localhost:~$
```

echo \$PATH

```
hduser@localhost:~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
hduser@localhost:~$ echo $PATH
/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/home/hduser/hadoop/bin:/home/hduser/hadoop/sbin:/home/hduser/.local/bin:/home/hduser/spark/bin
hduser@localhost:~$
```

JAVA WORD COUNT: ПРОЕКТ В IDEA

JAVA WORD COUNT: ЗАВИСИМОСТИ

Добавим в build.gradle зависимость на Apache Spark:

SparkDF - build.gradle (SparkDF)

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF build.gradle
Project build.gradle (SparkDF) ...
SparkDF ~/IdeaProjects/SparkDF
> .gradle
> .idea
> gradle
> src
  build.gradle
  gradlew
  gradlew.bat
  settings.gradle
> External Libraries
> Scratches and Consoles

build.gradle (SparkDF) x
1 plugins {
2     id 'java'
3 }
4
5 group 'org.example'
6 version '1.0-SNAPSHOT'
7
8 repositories {
9     mavenCentral()
10 }
11
12 java.toolchain.languageVersion = JavaLanguageVersion.of(8)
13
14 ext.sprk = "2.4.8"
15
16 dependencies {
17     compileOnly "org.apache.spark:spark-core_2.11:$sprk",
18     "org.apache.spark:spark-sql_2.11:$sprk"
19 }
20
21 dependencies{}
```

Build: Sync

SparkDF: finished At 24.11.2021, 23:31 3 sec, 931 ms > Task :prepareKotlinBuildScriptModel UP-TO-DATE

BUILD SUCCESSFUL in 3s

Favorites Structure

TODO Problems Terminal Build Dependencies Event Log

Shared indexes are downloaded for Maven library in 7 sec, 459 ms (42,4 MB) (moments ago) 18:45 LF UTF-8 4 spaces

## JAVA WORD COUNT: КОД

Создадим WordCount.java:

SparkDF - WordCount.java [SparkDF.main]

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF src main java WordCount
Project build.gradle (SparkDF) x WordCount.java x
SparkDF ~/IdeaProjects/Sp
> .gradle
> .idea
> build
> gradle
> src
  main
    > java
    > resources
  test
  build.gradle
  gradlew
  gradlew.bat
  ppkm_dataset.csv
  settings.gradle
> External Libraries
> Scratches and Consoles

WordCount.java x
1 import org.apache.spark.sql.Dataset;
2 import org.apache.spark.sql.Row;
3 import org.apache.spark.sql.SparkSession;
4
5 public class WordCount {
6
7     public static void main(String[] args) {
8         final String input = "ppkm_dataset.csv"; //args[0];
9         //final String output = args[1];
10
11         SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();
12
13         Dataset<Row> df = spark.read().option("header", "true").csv(input);
14
15         df.show();
16
17         spark.stop();
18     }
19 }
```

Run: SparkDF [WordCount.main()]

- SparkDF [WordCount.main()]: successful At 25.11.2021, 8 sec, 423 ms**
- WordCount.main() 5 warnings**
  - An illegal reflective access operation has occurred
  - Illegal reflective access by org.apache.spark.unsafe.Platform
  - Please consider reporting this to the maintainers of org.apache.spark
  - Use --illegal-access=warn to enable warnings of further illegal reflective access operations
  - All illegal access operations will be denied in a future release

24/11/25 00:37:54 INFO ShutdownHookManager: Shutdown hook called  
24/11/25 00:37:54 INFO ShutdownHookManager: Deleting directory /tmp/spark-9e3db9d2-1

class	comment
Ipositif	Kami siap laksana...
Ipositif	Siap melaksanakan...
Ipositif	Siap dukung dan s...
Ipositif	Langkah 3M ini su...
Ipositif	Siap amankan selu...
Ipositif	Siap utk di sosia...
Ipositif	Mendukung kebijak...

Run TODO Problems Terminal Build Dependencies Event Log

Rainbow CSV: You can edit Rainbow CSV settings in Settings > Editor > General > Rainbow CSV (2 minutes ago) 20:1 LF UTF-8 4 spaces

SparkDF - WordCount.java [SparkDF.main]

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF > src > main > java > WordCount > main
Project build.gradle (SparkDF) WordCount.java
SparkDF ~/IdeaProjects/Sp
  > .gradle
  > .idea
  > build
  > gradle
  > src
    > main
      > java
      > resources
    > test
    build.gradle
    gradlew
    gradlew.bat
    ppkm_dataset.csv
    settings.gradle
  > External Libraries
  > Scratches and Consoles
  > External Libraries
  > Scratches and Consoles

  2 import org.apache.spark.sql.Row;
  3   import org.apache.spark.sql.SparkSession;
  4
  5 public class WordCount {
  6
  7   public static void main(String[] args) {
  8     final String input = "ppkm_dataset.csv"; //args[0];
  9     //final String output = args[1];
 10
 11     SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();
 12
 13     Dataset<Row> df = spark.read().option("header", "true").csv(input);
 14
 15     df.printSchema();
 16
 17     spark.stop();
 18   }
 19 }

Run: SparkDF [WordCount.main()]
  ▾ ▲ SparkDF [WordCount.main()]: successful At 25.11. 6 sec, 769 ms
  ▾ ▲ :WordCount.main() 5 warnings 4 sec, 874 ms
    ▲ An illegal reflective access operation has occurred
    ▲ Illegal reflective access by org.apache.spark.unsafe.Platform
    ▲ Please consider reporting this to the maintainers of org.apache.spark
    ▲ Use --illegal-access=warn to enable warnings of further illegal reflective access operations
    ▲ All illegal access operations will be denied in a future release

  21/11/25 00:45:56 INFO ShutdownHookManager: Shutdown hook called
  21/11/25 00:45:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-eeac7f3e-3
  root
  |-- class: string (nullable = true)
  |-- comment: string (nullable = true)

  Deprecated Gradle features were used in this build, making it incompatible with Gradle 8.0.
  You can use '--warning-mode all' to show the individual deprecation warnings and details.

  See https://docs.gradle.org/7.1/userguide/command\_line\_interface.html#sec-command-line-deprecation

Run TODO Problems Terminal Build Dependencies Event Log
Rainbow CSV: You can edit Rainbow CSV settings in Settings > Editor > General > Rainbow CSV (8 minutes ago) 15:26 LF UTF-8 4 spaces

```

SparkDF - WordCount.java [SparkDF.main]

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF > src > main > java > WordCount > main
Project build.gradle (SparkDF) WordCount.java
SparkDF ~/IdeaProjects/Sp
  > .gradle
  > .idea
  > build
  > gradle
  > src
    > main
      > java
      > resources
    > test
    build.gradle
    gradlew
    gradlew.bat
    ppkm_dataset.csv
    settings.gradle
  > External Libraries
  > Scratches and Consoles
  > External Libraries
  > Scratches and Consoles

  1 import org.apache.spark.sql.Dataset;
  2 import org.apache.spark.sql.Row;
  3 import org.apache.spark.sql.SparkSession;
  4 import static org.apache.spark.sql.functions.*;
  5
  6 public class WordCount {
  7
  8   public static void main(String[] args) {
  9     final String input = "ppkm_dataset.csv"; //args[0];
 10     //final String output = args[1];
 11
 12     SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();
 13
 14     Dataset<Row> df = spark.read().option("header", "true").csv(input);
 15
 16     df.select(concat_ws( sep: " ", col( colName: "class"), col( colName: "comment"))).as("docs")
 17       .show();
 18
 19     spark.stop();
 20   }
 21 }

Run: SparkDF [WordCount.main()]
  ▾ ▲ SparkDF [WordCount.main()]: successful At 25.11. 6 sec, 923 ms
  ▾ ▲ :WordCount.main() 5 warnings 5 sec, 182 ms
    ▲ An illegal reflective access operation has occurred
    ▲ Illegal reflective access by org.apache.spark.unsafe.Platform
    ▲ Please consider reporting this to the maintainers of org.apache.spark
    ▲ Use --illegal-access=warn to enable warnings of further illegal reflective access operations
    ▲ All illegal access operations will be denied in a future release

  21/11/25 00:55:41 INFO MemoryStore: MemoryStore cleared
  21/11/25 00:55:41 INFO BlockManager: BlockManager stopped
  21/11/25 00:55:41 INFO BlockManagerMaster: BlockManagerMaster stopped
  21/11/25 00:55:41 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped
  +-----+
  |positif Kami siap...|
  |positif Siap mela...|
  |positif Siap duku...|
  |positif Langkah 3...|
  |positif Sian aman...|
  +-----+

Run TODO Problems Terminal Build Dependencies Event Log
Rainbow CSV: You can edit Rainbow CSV settings in Settings > Editor > General > Rainbow CSV (19 minutes ago) 17:33 LF UTF-8 4 spaces

```

SparkDF - WordCount.java [SparkDF.main]

```

import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.SparkSession;
import static org.apache.spark.sql.functions.*;

public class WordCount {
    public static void main(String[] args) {
        final String input = "ppkm_dataset.csv"; //args[0];
        //final String output = args[1];

        SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();

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

        df.select(concat_ws( sep: " ", col( colName: "class"), col( colName: "comment")).as("docs"))
            .select(split(col( colName: "docs"), pattern: "\\s").as("words"))
            .show();
    }
}

```

Run: SparkDF [:WordCount.main()]

- SparkDF [WordCount.main()]: successful At 25.11. 6 sec, 893 ms
- :WordCount.main() 5 warnings 5 sec, 187 ms
  - An illegal reflective access operation has occurred
  - Illegal reflective access by org.apache.spark.unsafe.Platform
  - Please consider reporting this to the maintainers of org.apache.spark
  - Use --illegal-access=warn to enable warnings of further illegal reflective access operations
  - All illegal access operations will be denied in a future release

```

21/11/25 00:58:21 INFO DAGScheduler: Job 1 finished: show at WordCount.java:18, took
21/11/25 00:58:21 INFO SparkUI: Stopped Spark web UI at http://192.168.1.68:4041
+-----+
|          words|
+-----+
|[positif, Kami, s...]
|[positif, Siap, m...]
|[positif, Siap, d...]
|[positif, Langkah...]
|[positif, Siap, a...]
|[positif, Siap, u...]
|[positif, Menduk...]

```

Event Log: 18:21 LF UTF-8 4 spaces

SparkDF - WordCount.java [SparkDF.main]

```

import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.SparkSession;
import static org.apache.spark.sql.functions.*;

public class WordCount {
    public static void main(String[] args) {
        final String input = "ppkm_dataset.csv"; //args[0];
        //final String output = args[1];

        SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();

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

        df.select(concat_ws( sep: " ", col( colName: "class"), col( colName: "comment")).as("docs"))
            .select(split(col( colName: "docs"), pattern: "\\s").as("words"))
            .printSchema();
    }
}

```

Run: SparkDF [:WordCount.main()]

- SparkDF [WordCount.main()]: successful At 25.11. 6 sec, 727 ms
- :WordCount.main() 5 warnings 5 sec, 26 ms
  - An illegal reflective access operation has occurred
  - Illegal reflective access by org.apache.spark.unsafe.Platform
  - Please consider reporting this to the maintainers of org.apache.spark
  - Use --illegal-access=warn to enable warnings of further illegal reflective access operations
  - All illegal access operations will be denied in a future release

```

21/11/25 00:59:18 INFO MemoryStore: Block broadcast_2 stored as values in memory les
21/11/25 00:59:18 INFO MemoryStore: Block broadcast_2_piece0 stored as bytes in memo
21/11/25 00:59:18 INFO BlockManagerInfo: Added broadcast_2_piece0 in memory on 192.1
21/11/25 00:59:18 INFO SparkContext: Created broadcast 2 from csv at WordCount.java:
21/11/25 00:59:18 INFO FileSourceScanExec: Planning scan with bin packing, max size:
root
  |-- words: array (nullable = false)
  |    |-- element: string (containsNull = true)

21/11/25 00:59:18 INFO SparkUI: Stopped Spark web UI at http://192.168.1.68:4041
21/11/25 00:59:18 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint
21/11/25 00:59:18 INFO MemoryStore: MemoryStore cleared

```

Event Log: 18:36 LF UTF-8 4 spaces

SparkDF - WordCount.java [SparkDF.main]

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF > src > main > java > WordCount > main
Project build.gradle (SparkDF) > WordCount.java
1 import org.apache.spark.sql.Dataset;
2 import org.apache.spark.sql.Row;
3 import org.apache.spark.sql.SparkSession;
4 import static org.apache.spark.sql.functions.*;
5
6 public class WordCount {
7
8     public static void main(String[] args) {
9         final String input = "ppkm_dataset.csv"; //args[0];
10        //final String output = args[1];
11
12        SparkSession spark = SparkSession.builder().master("local[*]").getOrCreate();
13
14        Dataset<Row> df = spark.read().option("header", "true").csv(input);
15
16        df.select(concat_ws(" ", col("class"), col("comment")).as("docs"))
17            .select(split(col("docs"), pattern("\\s+").as("words")))
18            .select(explode(col("words")).as("word"))
19            .show();
20
21        spark.stop();
22    }
23}

```

Run: SparkDF [:WordCount.main()]

- SparkDF [WordCount.main()]: successful At 25.11.6 sec, 925 ms
  - :WordCount.main() 5 warnings
    - An illegal reflective access operation has occurred
    - Illegal reflective access by org.apache.spark.unsafe.Platform
    - Please consider reporting this to the maintainers of org.apache.spark
    - Use --illegal-access=warn to enable warnings of further illegal reflective access operations
    - All illegal access operations will be denied in a future release

Structure Favorites Event Log

Rainbow CSV: You can edit Rainbow CSV settings in Settings > Editor > General > Rainbow CSV (24 minutes ago)

19:29 LF UTF-8 4 spaces

Соберем jar:

SparkDF - WordCount.java [SparkDF.main]

```

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help
SparkDF > build > libs > SparkDF-1.0-SNAPSHOT.jar
Project build.gradle (SparkDF) > WordCount.java
1 import org.apache.spark.sql.Row;
2 import org.apache.spark.sql.SaveMode;
3 import org.apache.spark.sql.SparkSession;
4 import static org.apache.spark.sql.functions.*;
5
6 public class WordCount {
7
8     public static void main(String[] args) {
9         final String input = args[0];
10        final String output = args[1];
11
12        SparkSession spark = SparkSession.builder().getOrCreate();
13
14        Dataset<Row> df = spark.read().option("header", "true").csv(input);
15
16        Dataset<Row> wc = df.select(concat_ws(" ", col("class"),
17            .select(split(col("docs"), pattern("\\s+").as("words"))
18            .select(explode(col("words")).as("word"))
19            .groupBy("word").count());
20
21        wc.write().mode(SaveMode.Overwrite).csv(output);
22
23        spark.stop();
24}

```

Gradle

- SparkDF
  - Tasks
    - assemble
    - build
    - buildDepends
    - buildNeeded
    - classes
    - clean
    - jar
    - testClasses
    - build setup
    - documentation
    - help
    - other
    - verification
    - Dependencies

Run: SparkDF [build]

- SparkDF [build]: successful At 25.11.2021, 1:12 1 sec, 776 ms
  - > Task :test NO-SOURCE
  - > Task :check UP-TO-DATE
  - > Task :build

BUILD SUCCESSFUL in 1s  
2 actionable tasks: 2 executed  
1:12:12: Task execution finished 'build'.

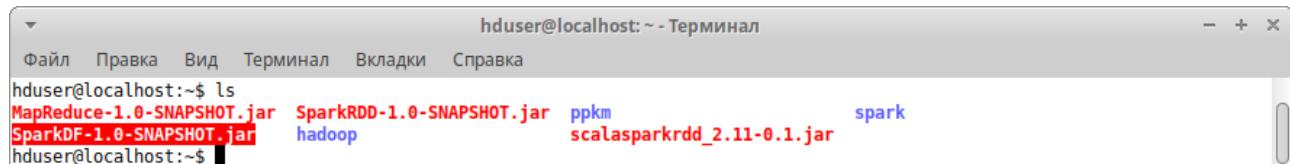
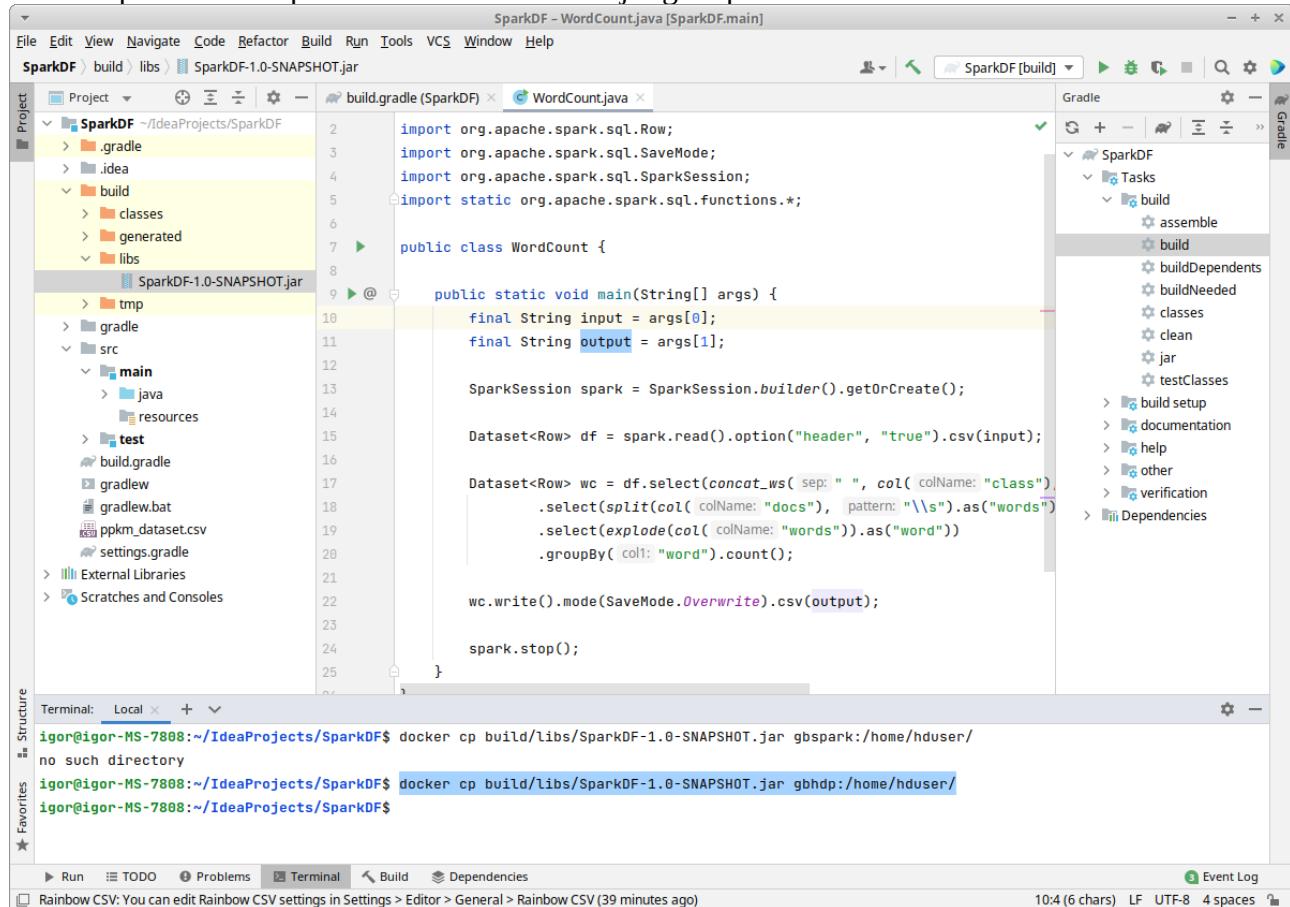
Structure Favorites Event Log

Rainbow CSV: You can edit Rainbow CSV settings in Settings > Editor > General > Rainbow CSV (36 minutes ago)

10:4 (6 chars) LF UTF-8 4 spaces

Перенесем его на кластер:

```
docker cp build/libs/SparkDF-1.0-SNAPSHOT.jar gbdhp:/home/hduser/
```



```
ls -la
```

```
hduser@localhost:~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
hduser@localhost:~$ ls -la
total 80
drwxr-xr-x 1 hduser hduser 4096 Nov 24 22:15 .
drwxr-xr-x 1 root   root   4096 Nov 22 20:40 ..
-rw----- 1 hduser hduser 1535 Nov 24 19:37 .bash_history
-rw-r--r-- 1 hduser hduser 220 Apr  4 2018 .bash_logout
-rw-r--r-- 1 hduser hduser 3889 Nov 24 17:56 .bashrc
drwx----- 2 hduser hduser 4096 Nov 22 20:43 .cache
-rw-r--r-- 1 hduser hduser 807 Apr  4 2018 .profile
drwx----- 2 hduser hduser 4096 Nov 22 20:40 .ssh
-rw-r--r-- 1 hduser hduser    0 Nov 22 20:40 .sudo_as_admin_successful
-rw----- 1 hduser hduser 1124 Nov 24 17:56 .viminfo
-rw-rw-r-- 1 hduser hduser 4441 Nov 22 22:32 MapReduce-1.0-SNAPSHOT.jar
-rw-rw-r-- 1 hduser hduser 1445 Nov 24 22:12 SparkDF-1.0-SNAPSHOT.jar
-rw-rw-r-- 1 hduser hduser 2380 Nov 24 18:35 SparkRDD-1.0-SNAPSHOT.jar
drwxr-xr-x 1 hduser hduser 4096 Nov 22 20:43 hadoop
drwxr-xr-x 2 hduser hduser 4096 Jan  1 1970 ppkm
-rw-rw-r-- 1 hduser hduser 6604 Nov 24 19:09 scalasparkrdd_2.11-0.1.jar
drwxr-xr-x 13 hduser hduser 4096 May  8 2021 spark
hduser@localhost:~$
```

Запускаем приложение:

```
$ spark-submit \
--class WordCount \
--master yarn \
--deploy-mode cluster \
SparkDF-1.0-SNAPSHOT.jar \
/usr/hduser/ppkm/ppkm_dataset.csv /user/hduser/ppkm-df-java
```



Application application\_1637785024774\_0001 — Mozilla Firefox

Файл Правка Вид Журнал Закладки Инструменты Справка

Application application\_1637785024774\_0001 +

localhost:8088/cluster/app/application\_1637785024774\_0001

Logged in as: dr.who

**Application**  
**application\_1637785024774\_0001**

**Cluster**

- About
- Nodes
- Node Labels
- Applications**
  - NEW
  - NEW\_SAVING
  - SUBMITTED
  - ACCEPTED
  - RUNNING
  - FINISHED
  - FAILED
  - KILLED
- Scheduler

**Tools**

**Application Overview**

User: [hduser](#)  
Name: WordCount  
Application Type: SPARK  
Application Tags:  
Application Priority: 0 (Higher Integer value indicates higher priority)  
YarnApplicationState: FINISHED  
Queue: default  
FinalStatus Reported by AM: SUCCEEDED  
Started: Wed Nov 24 22:21:15 +0000 2021  
Launched: Wed Nov 24 22:21:16 +0000 2021  
Finished: Wed Nov 24 22:22:39 +0000 2021  
Elapsed: 1mins, 24sec  
Tracking URL: History  
Log Aggregation Status: DISABLED  
Application Timeout (Remaining Time): Unlimited  
Diagnostics:  
Unmanaged Application: false  
Application Node Label expression: <Not set>  
AM container Node Label expression: <DEFAULT\_PARTITION>

**Application Metrics**

Total Resource Preempted: <memory:0, vCores:0>  
Total Number of Non-AM Containers Preempted: 0

Смотрим результат:

\$ hdfs dfs -ls /user/hduser/ppkm-df-java

hduser@localhost: ~ - Терминал

Файл Правка Вид Терминал Вкладки Справка

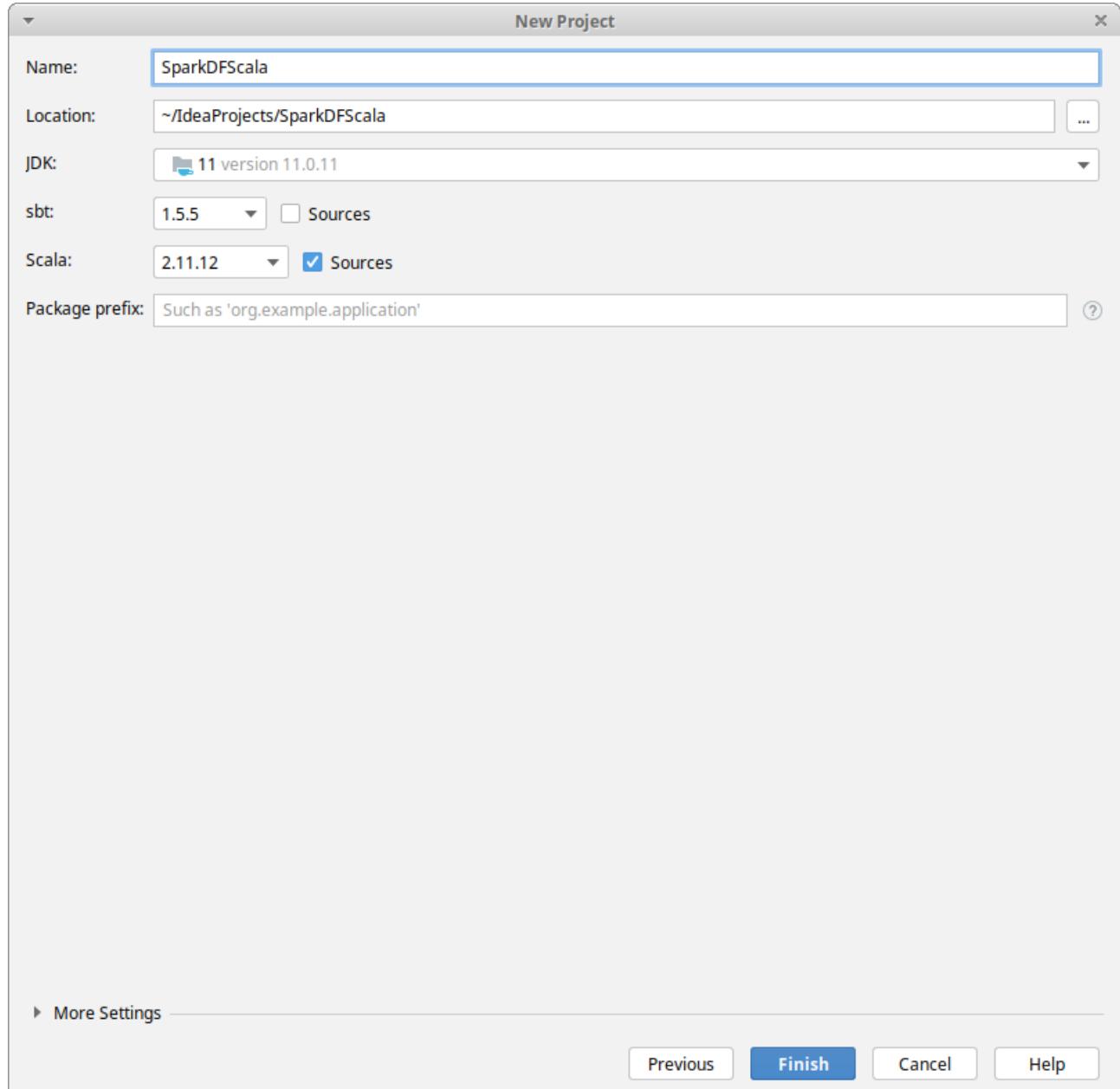
```
-rw-r--r-- 1 hduser supergroup 171 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00184-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 108 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00185-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 93 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00186-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 89 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00187-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 222 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00188-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 96 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00189-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 141 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00190-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 114 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00191-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 109 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00192-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 122 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00193-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 130 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00194-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 102 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00195-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 139 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00196-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 172 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00197-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 105 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00198-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
-rw-r--r-- 1 hduser supergroup 137 2021-11-24 22:22 /user/hduser/ppkm-df-java/part-00199-7d5d387d-af2a-4b0b-bbbe-4fe7ab544b20-c000.csv
```

hduser@localhost:~\$

\$ hdfs dfs -cat /user/hduser/ppkm-df-java/\*

hduser@localhost: ~ Терминал  
Файл Правка Вид Терминал Вкладки Справка  
penghasilannya,1  
mata,1  
#pelanggarpsbb,1  
lindungi,1  
menanggapi,2  
Kakak...Aamiin,1  
Masha,1  
Munir,1  
"7/INSTR/2021,",4  
lo,4  
dengar,1  
msih,1  
diganti,1  
angka,3  
Gencarkan,1  
bali.,1  
Targetkan,1  
PENANGANAN,1  
Wakapolda,1  
bekerumun,1  
setengah,1  
"didiemin,",1  
Bpk,1  
kali,1  
....yg,1  
<https://t.co/HwFdBrSvVI>,1  
kawasan,1  
jenazah,1  
tapi,7  
via,1  
mau,10  
beban,1  
goreng,1  
muter2,1  
masa,2  
nafkah,1  
"revisi,",1  
mencegah,9  
harian,1  
namun,2  
Tabarakallah,1  
ayo,1  
Gimana,2  
"pak,malah",1  
petugah,1  
massal...,1  
Digelar,1  
Rembang,1  
dibikin,1  
gara,13  
Weekend,1  
tertibkan,1  
luas,1  
stop,1  
yaitu,1  
he,2  
perekonomian,1  
<https://t.co/AMtVMfHKtl>,1  
Covid-19#indonesia,2  
Termasuk,2  
boleh,4  
tetpakasa,1  
pusat,1  
<https://t.co/s0hi56lRZA>,3  
kalo,5  
baiknya,1  
telat,1  
agar,1  
Kami,1  
#Dirumahaja,1  
Kegiatan,39  
Nomor:,4  
👉,1  
nyinyirnya,1  
merugikan,2  
"Dago,",1  
Malioboro,1  
sepeda,1  
kekuatan,1  
Insya,1  
dilaksanakan,2  
Perpanjangan,8  
"Mitra,",2  
#beritajabar,1  
Belum,1  
mentok,1  
kebudayaan,1  
DP,3  
empat,1  
Percuma,1  
maksud,1  
hduser@localhost:~\$

## SCALA WORD COUNT: ПРОЕКТ В IDEA



## SCALA WORD COUNT: ЗАВИСИМОСТИ

Добавим в build.sbt зависимость на Apache Spark:

```
name := "SparkDFScala"
version := "0.1"
scalaVersion := "2.11.12"

javacOptions ++= Seq("-source", "1.8", "-target", "1.8")
val sparkVersion = "2.4.8"
libraryDependencies ++= Seq(
  "org.apache.spark" %% "spark-core" % sparkVersion,
  "org.apache.spark" %% "spark-sql" % sparkVersion)
```

SparkDFScala - WordCount.scala [SparkDFScala]

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help

SparkDFScala > src > main > scala > WordCount.scala

Project

SparkDFScala ~/IdeaProjects/SparkDFScala

.bsp .idea

project [SparkDFScala-build] sources root

src

main

scala

WordCount

test

target

build.sbt

External Libraries

Scratches and Consoles

build.sbt

WordCount.scala

```
1 import org.apache.spark.sql.SparkSession
2 import org.apache.spark.sql.functions._
3
4 object WordCount {
5   def main(args: Array[String]): Unit = {
6     val input = args(0)
7     val output = args(1)
8
9     val spark = SparkSession.builder().getOrCreate()
10    import spark.sqlContext.implicits._
11
12    spark.read.option("header", "true").csv(input)
13      .select(concat_ws(" ", exprs = $"class", $"comment") as "docs")
14      .select(split(str = $"docs", pattern = "\\s") as "words")
15      .select(explode($"words") as "word")
16      .groupBy(cols = $"word").count()
17      .write.mode(saveMode = "overwrite").csv(output)
18
19    spark.stop()
20  }
21}
```

WordCount > main(args: Array[String])

TODO Problems Terminal sbt shell Build Dependencies Event Log

18:1 LF UTF-8 2 spaces

SparkDFScala - WordCount.scala [SparkDFScala]

File Edit View Navigate Code Refactor Build Run Tools VCS Window Help

SparkDFScala > src > main > scala > WordCount.scala

Project

SparkDFScala ~/IdeaProjects/SparkDFScala

.bsp .idea

project [SparkDFScala-build] sources root

src

main

scala

WordCount

test

target

build.sbt

External Libraries

Scratches and Consoles

build.sbt

WordCount.scala

sbt

+ makeIvyXmlConfiguration  
makeIvyXmlLocalConfiguration  
makePom  
managedClasspath  
managedResources  
managedSources  
manipulateBytecode  
mappings  
moduleSettings  
otherResolvers  
package  
packageBin  
packageCache  
packageConfiguration  
packagedArtifact  
packagedArtifacts  
packageDoc  
packageOptions  
packageSrc  
pickleProducts  
previousCompile  
printWarnings  
productDirectories  
products  
projectDependencies  
projectDescriptors  
projectResolver  
publish  
publishConfiguration  
publisher  
publishLocal  
publishLocalConfiguration

```
1 import org.apache.spark.sql.SparkSession
2 import org.apache.spark.sql.functions._
3
4 object WordCount {
5   def main(args: Array[String]): Unit = {
6     val input = args(0)
7     val output = args(1)
8
9     val spark = SparkSession.builder().getOrCreate()
10    import spark.sqlContext.implicits._
11
12    spark.read.option("header", "true").csv(input)
13      .select(concat_ws(" ", exprs = $"class")
14      .select(split(str = $"docs", pattern = "\\s"))
15      .select(explode($"words") as "word")
16      .groupBy(cols = $"word").count()
17      .write.mode(saveMode = "overwrite").csv(output)
18
19    //wc.write.
20
21    spark.stop()
22  }
23}
```

WordCount > main(args: Array[String])

TODO Problems Terminal sbt shell Build Dependencies Event Log

19:7 LF UTF-8 2 spaces

The screenshot shows the IntelliJ IDEA interface with the following details:

- Project View:** Shows the project structure under "SparkDFScala". The "target" directory contains "scala-2.11" which includes "sparkdfscala\_2.11-0.1.jar".
- Code Editor:** Displays the content of `WordCount.scala`. The code reads a CSV file with columns "class" and "comment", splits it into words, counts them, and writes the results back to a CSV file.
- Sbt Shell:** Shows the command-line output of the sbt build process:

```
[info] set current project to SparkDFScala (in build file:/home/igor/IdeaProjects/SparkDFScala/)  
[IJ]{file:/home/igor/IdeaProjects/SparkDFScala/}sparkdfscala/package  
[info] compiling 1 Scala source to /home/igor/IdeaProjects/SparkDFScala/target/scala-2.11/classes ...  
[info] done compiling  
[success] Total time: 5 s, completed 25 Mar 2021, 1:44:53  
[IJ]  
>
```
- Bottom Bar:** Includes tabs for "TODO", "Problems", "Terminal", "sbt shell", "Build", and "Dependencies".

docker cp target/scala-2.11/sparkdfscala\_2.11-0.1.jar gbdhp:/home/hduser/

The screenshot shows the IntelliJ IDEA interface with the following details:

- Project View:** Shows the project structure under "SparkDFScala". The "target" directory contains "scala-2.11" which includes "sparkdfscala\_2.11-0.1.jar".
- Code Editor:** Displays the content of `WordCount.scala`, identical to the first screenshot.
- Terminal:** Shows the command-line output of the "docker cp" command:

```
igor@igor-MS-7808:~/IdeaProjects/SparkDFScala$ docker cp target/scala-2.11/sparkdfscala_2.11-0.1.jar gbdhp:/home/hduser/  
igor@igor-MS-7808:~/IdeaProjects/SparkDFScala$
```
- Bottom Bar:** Includes tabs for "TODO", "Problems", "Terminal", "sbt shell", "Build", and "Dependencies".

```

hduser@localhost:~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
hduser@localhost:~$ ls
MapReduce-1.0-SNAPSHOT.jar  SparkRDD-1.0-SNAPSHOT.jar  ppkm
SparkDF-1.0-SNAPSHOT.jar  hadoop  scalasparkrdd_2.11-0.1.jar  spark
hduser@localhost:~$ █

```

## SCALA WORD COUNT: ЗАПУСК

Запускаем приложение:

```

spark-submit \
--class WordCount \
--master yarn \
--deploy-mode cluster \
sparkdfscala_2.11-0.1.jar \
/usr/hduser/ppkm/ppkm_dataset.csv /user/hduser/ppkm-df-scala

```

The screenshot shows the Apache Spark 2.4.8 application UI running in a Mozilla Firefox browser. The URL is `localhost:8088/proxy/application_1637785024774_0003/`. The UI has tabs for Jobs, Stages, Storage, Environment, Executors, and SQL. The Jobs tab is selected, displaying the "Spark Jobs" section. It shows the following information:

- User: hduser
- Total Uptime: 56 s
- Scheduling Mode: FIFO
- Active Jobs: 1
- Completed Jobs: 1

Under the "Active Jobs (1)" section, there is one entry:

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
1	csv at WordCount.scala:17 csv at WordCount.scala:17 (kill)	2021/11/24 22:58:51	20 s	1/2	108/201 (2 running)

Under the "Completed Jobs (1)" section, there is one entry:

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	csv at WordCount.scala:12 csv at WordCount.scala:12	2021/11/24 22:58:49	2 s	1/1	1/1





```
hduser@localhost: ~ - Терминал
Файл Правка Вид Терминал Вкладки Справка
https://t.co/AMtVMfHKTl,1
Covid-19#indonesia,2
Termasuk,2
boleh,4
tetpakasa,1
pusat,1
https://t.co/s0hi56lRZA,3
kalo,5
baiknya,1
telat,1
agar,1
Kami,1
#dirumahaja,1
Kegiatan,39
Nomor:,4
👉,1
nyinyirnya,1
merugikan,2
"Dago,",1
Malioboro,1
sepeda,1
kekuatan,1
Insya,1
dilaksanakan,2
Perpanjangan,8
"Mitra,",2
#beritajabar,1
Belum,1
mentok,1
kebudayaan,1
DP,3
empat,1
Percuma,1
maksud,1
hduser@localhost:~$ █
```