



SAKARYA
ÜNİVERSİTESİ

BIG DATA

TOO BIG TO IGNORE

SÜMEYYE KAYNAK

OUTLINE

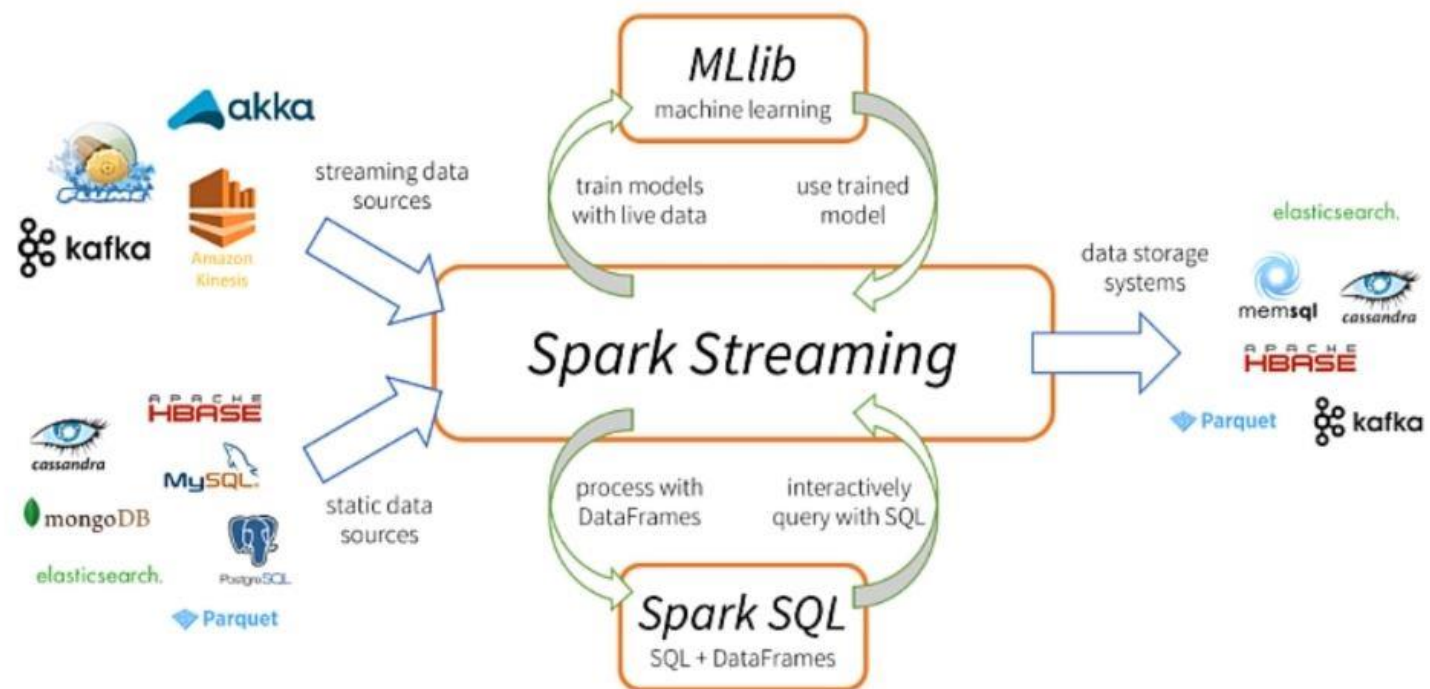


Spark

SPARK

- It is a big data library developed with Scala to analyze big data with distributed processing. Spark works in-memory.
- Therefore, there is no storage unit. It performs data analysis on RAM.

SPARK



SPARK AND HADOOP

Storage:

Hadoop=> HDFS

Spark=> There is no internal storage. But thanks to integration, data can be stored in different technologies.

Speed:

Spark is 100 times faster than Hadoop

Difficult:

Spark => It is easy to program thanks to Spark RDD and SQL.

Hadoop => Developing Map-reduce is difficult.

SPARK AND HADOOP

Management:

Hadoop=> Yarn

Spark=> provides its own management

Real-time analysis:

Hadoop=> There is no real-time analysis tool.

Spark => With Spark streaming, millions of data per second can be analyzed instantly.

SPARK MLLIB

- has an extra library (Spark MLlib) where you can apply machine learning techniques in Spark.
- Estimation and classification can be made with techniques such as logistic regression, K Means, K NN, Naive Bayes.

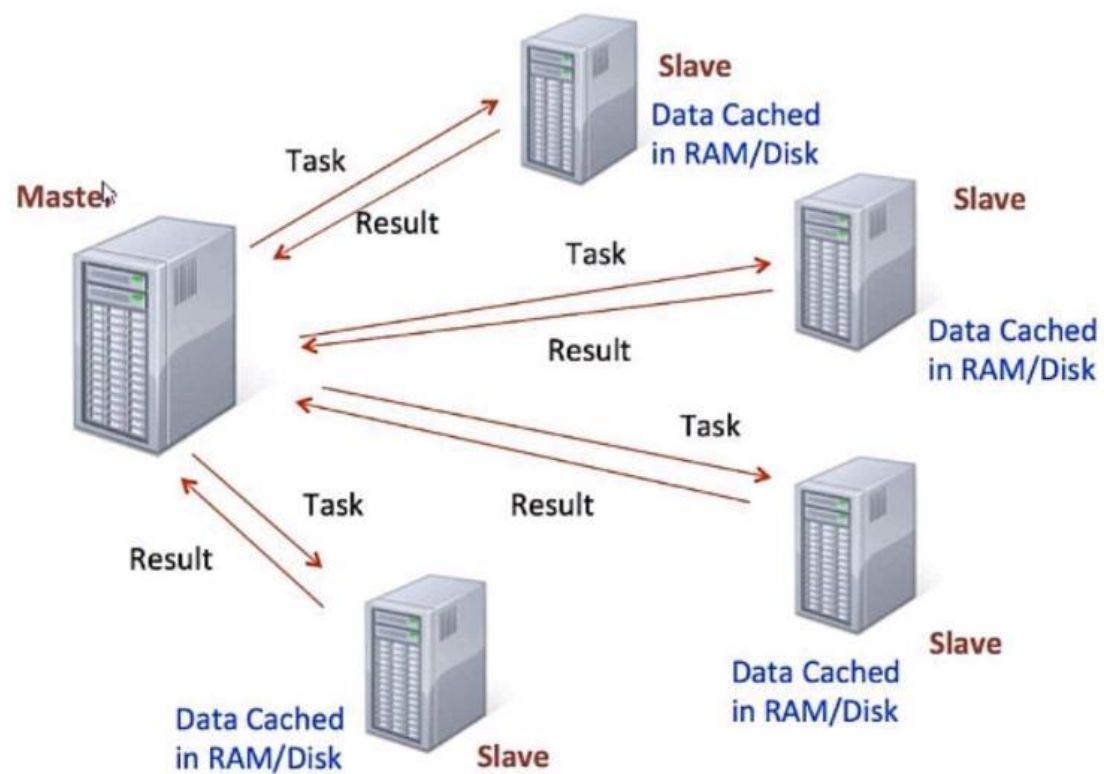
SPARK STREAMING

- Instant, real-time data analysis is performed with Spark's streaming library.

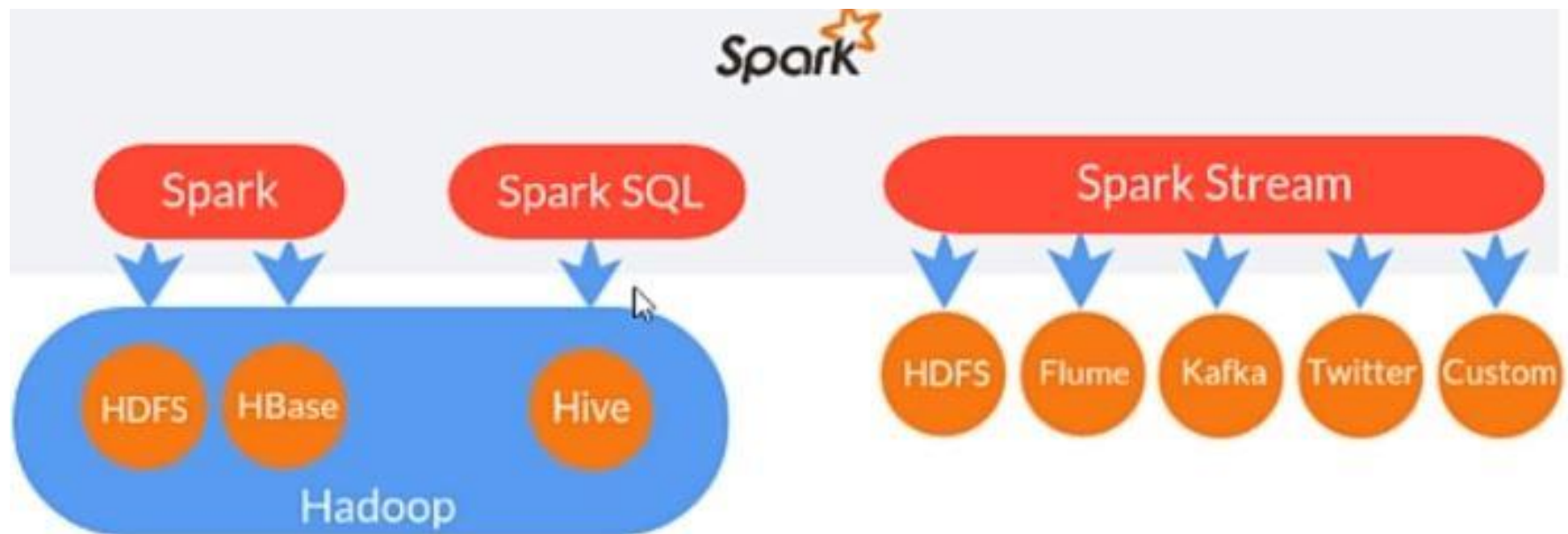
SPARK SQL

- Spark has a library for performing big data analysis with SQL-based queries.

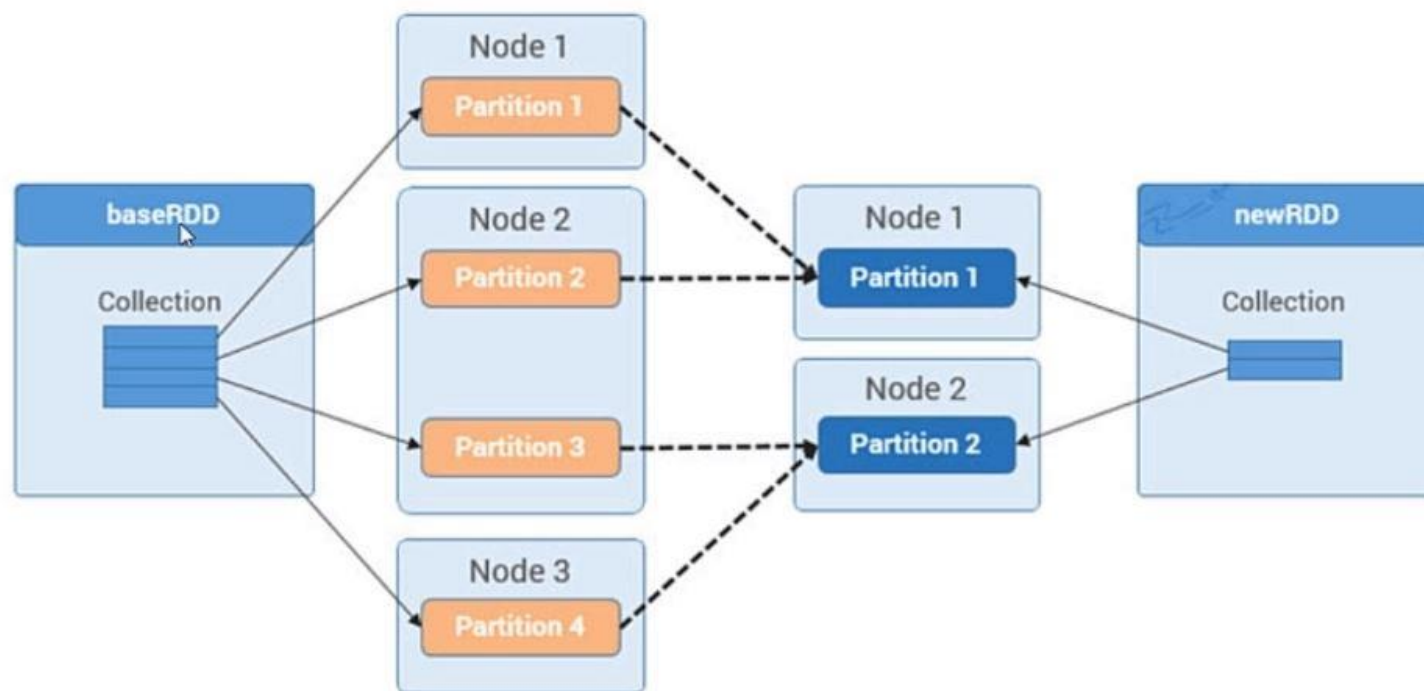
SPARK ARCHITECTURE



SPARK



SPARK RDD



SPARK RDD

[1,2,.....,99]

List

1,2,...

1,2,...

26,28,...

34,35,...

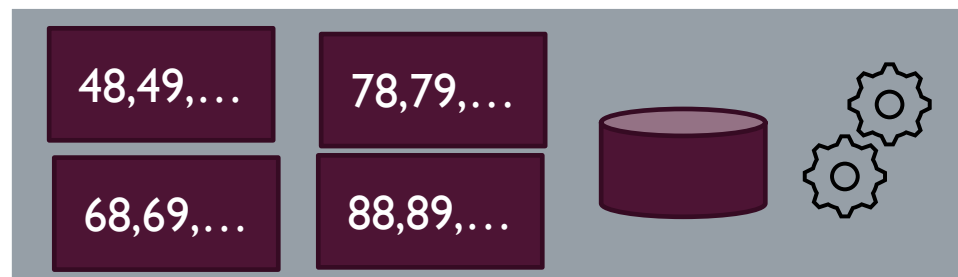
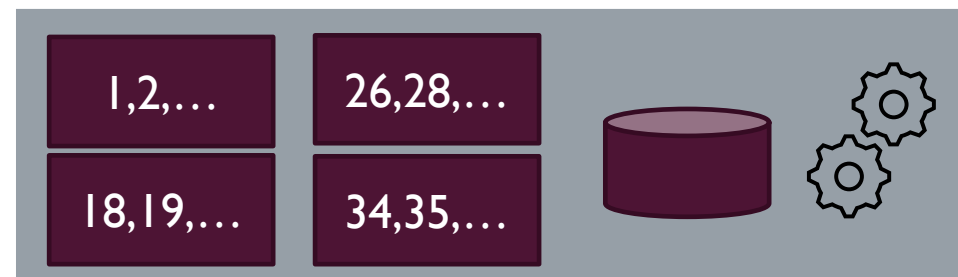
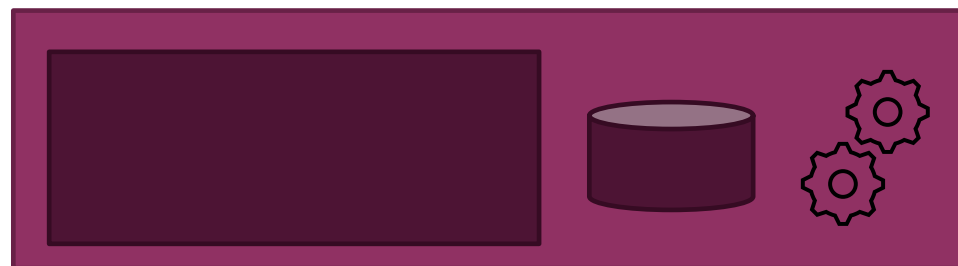
48,49,...

68,69,...

78,79,...

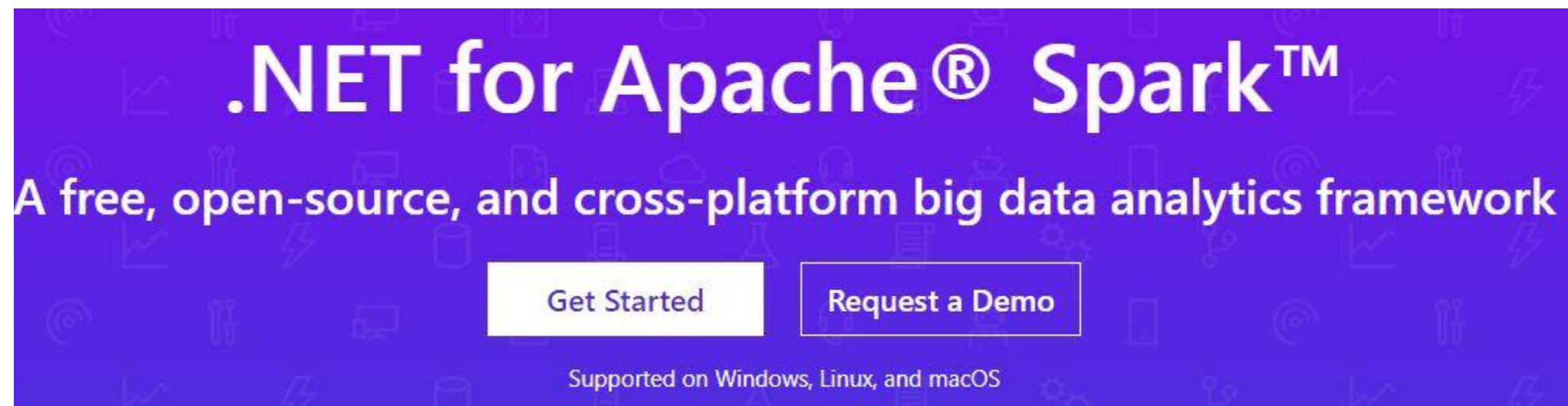
88,89,...

RDD



Physical

.NET FOR APACHE SPARK APP

A banner with a blue background featuring a repeating pattern of white icons related to data and technology. The text is centered and reads: ".NET for Apache® Spark™" in a large, bold, white font. Below this, in a smaller white font, is "A free, open-source, and cross-platform big data analytics framework". At the bottom, there are two white buttons with black text: "Get Started" and "Request a Demo". Below the buttons, in a small white font, is "Supported on Windows, Linux, and macOS".

.NET for Apache® Spark™

A free, open-source, and cross-platform big data analytics framework

[Get Started](#) [Request a Demo](#)

Supported on Windows, Linux, and macOS

Install .NET

To start building .NET apps, download and install the .NET SDK (Software Development Kit).

[Download .NET 3.1 SDK \(64-bit\)](#)

[32-bit download](#)

.NET FOR APACHE SPARK APP

```
C:\Users\Sumeyye>dotnet

Usage: dotnet [options]
Usage: dotnet [path-to-application]

Options:
  -h|--help          Display help.
  --info             Display .NET information.
  --list-sdks        Display the installed SDKs.
  --list-runtimes    Display the installed runtimes.

path-to-application:
  The path to an application .dll file to execute.
```

.NET FOR APACHE SPARK APP-INSTALL JAVA

1. Visit [Java SE Development Kit 8 Downloads](#).
2. Click the download for **Windows x64**.
3. Review the License agreement and accept it if you agree.
4. Oracle requires an account to download the JDK. So, sign in to an existing account or complete the account registration process to start downloading.
5. Once the download completes, run the installer using the default settings.

Once you've installed, open a **new** command prompt and run the following command:

```
C:\Users\Sumeyye>java -version
java version "11.0.13" 2021-10-19 LTS
Java(TM) SE Runtime Environment 18.9 (build 11.0.13+10-LTS-370)
Java HotSpot(TM) 64-Bit Server VM 18.9 (build 11.0.13+10-LTS-370, mixed mode)
```


.NET FOR APACHE SPARK APP-INSTALL JAVA

To download Java 8, click [Download Java](#)

1. Click the download for **Windows x64**.
2. Review the License agreement and accept it if you agree.
3. Oracle requires an account to download the JDK. So, sign in to an existing account or complete the account registration process to start downloading.
4. Once the download completes, run the installer using the default settings.



.NET FOR APACHE SPARK APP-INSTALL JAVA JDK8

JDK 8 software is licensed under the [Oracle Technology Network License Agreement for Oracle Java SE](#)

[JDK 8u311 checksum](#)

Linux **macOS** **Solaris** **Windows**

Product/file description

x86 Installer

x64 Installer

Documentation Download


You must accept the [Oracle Technology Network License Agreement](#) for Oracle Java SE to download this software.



I reviewed and accept the Oracle Technology Network License Agreement for Oracle Java SE

Required

You will be redirected to the login screen in order to download the file.

Download [jdk-8u311-windows-x64.exe](#) 

36.exe

4.exe

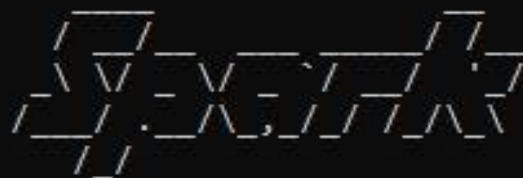
.NET FOR APACHE SPARK APP-DOWNLOAD APACHE SPARK

1. Apache Spark is downloaded as a .tgz file.
<https://archive.apache.org/dist/spark/spark-3.0.1/spark-3.0.1-bin-hadoop2.7.tgz>
2. Extract Apache Spark files in C directories
3. Open a **new** command prompt, run the following commands to set the environment variables used to locate Apache Spark:

```
C:\Users\Sumeyye>setx HADOOP_HOME C:\bin\spark-3.0.1-bin-hadoop2.7\  
SUCCESS: Specified value was saved.  
  
C:\Users\Sumeyye>setx SPARK_HOME C:\bin\spark-3.0.1-bin-hadoop2.7\  
SUCCESS: Specified value was saved.
```

CHECK SPARK VERSION

```
C:\bin\spark-3.0.1-bin-hadoop2.7\bin>spark-submit --version
```



version 3.0.1

```
Using Scala version 2.12.10, Java HotSpot(TM) 64-Bit Server VM, 11.0.13
Branch HEAD
Compiled by user ubuntu on 2020-08-28T07:36:48Z
Revision 2b147c4cd50da32fe2b4167f97c8142102a0510d
Url https://gitbox.apache.org/repos/asf/spark.git
Type --help for more information.
```

MICROSOFT.SPARK.WORKER

1. Download the Microsoft.Spark.Worker release from the .NET for Apache Spark GitHub repository:
(https://github.com/dotnet/spark/releases/download/v1.0.0/Microsoft.Spark.Worker.netcoreapp3.1.win-x64-1.0.0.zip?WT.mc_id=dotnet-35129-website)
2. Enter C:\bin in the Extract to field.

INSTALL WINUTILS

1. .NET for Apache Spark requires WinUtils to be installed alongside Apache Spark.

https://github.com/steveloughran/winutils/raw/master/hadoop-2.7.1/bin/winutils.exe?WT.mc_id=dotnet-35129-website

Once **winutils.exe** downloads, copy it into **C:\bin\spark-3.0.1-bin-hadoop2.7\bin**.

SET DOTNET_WORKER_DIR

- This is used by .NET apps to locate .NET for Apache Spark.

```
C:\bin\spark-3.0.1-bin-hadoop2.7\bin>setx DOTNET_WORKER_DIR "C:\bin\Microsoft.Spark.Worker-1.0.0"  
SUCCESS: Specified value was saved.
```

CREATE YOUR APP

```
> dotnet new console -f netcoreapp3.1 -o mySparkApp
```

```
> cd mySparkApp
```



Open visual studio and create a console application.


INSTALL THE NUGET PACKAGE

NuGet: mySparkApp Program.cs


[Browse](#) Installed Updates

spark ☐ Include prerelease Package source: nuget.org

	Spark by Louis DeJardin & Robert Greyling, 520K downloads Spark View Engine Core.	v1.8.1
	Microsoft.Spark by Microsoft, 361K downloads .NET for Apache Spark	v2.0.0

Microsoft.Spark  [nuget.org](#)

Version: Latest stable 2.0.0

 Options

RUNNING SPARK

```
Select Windows PowerShell
PS C:\Users\Sumeyye\source\repos\mySparkApp\mySparkApp> spark-submit --class org.apache.spark.deploy.dotnet.DotnetRunner
--master local .\bin\Debug\netcoreapp3.1\microsoft-spark-3-0_2.12-2.0.0.jar debug
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.spark.unsafe.Platform (file:/C:/bin/spark-3.0.1-bin-hadoop2.7/jars/spark-unsafe_2.12-3.0.1.jar) to constructor java.nio.DirectByteBuffer(long,int)
WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
21/12/12 15:34:40 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
21/12/12 15:34:40 INFO DotnetRunner: Starting DotnetBackend with .
21/12/12 15:34:40 INFO DotnetBackend: The number of DotnetBackend threads is set to 10.
21/12/12 15:34:42 INFO DotnetRunner: Port number used by DotnetBackend is 5567
*****
* .NET Backend running debug mode. Press enter to exit *
*****
```

ADD DATA FILE

Create an `input.txt` file in your `mySparkApp` directory, containing the following text.

```
input.txt
```

```
Hello World
```

```
This .NET app uses .NET for Apache Spark
```

```
This .NET app counts words with Apache Spark
```

CODE YOUR APP

```
// Create a Spark session
SparkSession spark = SparkSession
    .Builder()
    .AppName("word_count_sample")
    .GetOrCreate();

// Create initial DataFrame
DataFrame dataframe = spark.Read().Text("input.txt");

// Count words
DataFrame words = dataframe
    .Select(Functions.Split(Functions.Col("value"), " ").Alias("words"))
    .Select(Functions.Explode(Functions.Col("words")))
    .Alias("word"))
    .GroupBy("word")
    .Count()
    .OrderBy(Functions.Col("count").Desc());

// Show results
words.Show();

// Stop Spark session
spark.Stop();
```

RUNNING APP

```
Microsoft Visual Studio Debug Console
[2021-12-12T12:34:55.8163166Z] [DESKTOP-GK3VE45] [Info] [ConfigurationService] 'DOTNETBACKEND_PORT' environment variable
is not set.
[2021-12-12T12:34:55.8415011Z] [DESKTOP-GK3VE45] [Info] [ConfigurationService] Using port 5567 for connection.
[2021-12-12T12:34:55.8639466Z] [DESKTOP-GK3VE45] [Info] [JvmBridge] JvmBridge port is 5567
[2021-12-12T12:34:55.8806514Z] [DESKTOP-GK3VE45] [Info] [JvmBridge] The number of JVM backend thread is set to 10. The m
ax number of concurrent sockets in JvmBridge is set to 7.
+-----+-----+
| word|count|
+-----+-----+
|.NET| 3|
|Apache| 2|
| app| 2|
| This| 2|
| Spark| 2|
| World| 1|
| counts| 1|
| for| 1|
| words| 1|
| with| 1|
| Hello| 1|
| uses| 1|
+-----+-----+

C:\Users\Sumeyye\source\repos\mySparkApp\mySparkApp\bin\Debug\netcoreapp3.1\mySparkApp.exe (process 19780) exited with c
ode 0.
To automatically close the console when debugging stops, enable Tools->Options->Debugging->Automatically close the conso
le when debugging stops.
Press any key to close this window . . .
```

DATA LOADING STAGES

Installing from local computer

```
DataFrame wordRDD = spark.Read().Text("input.txt");
```

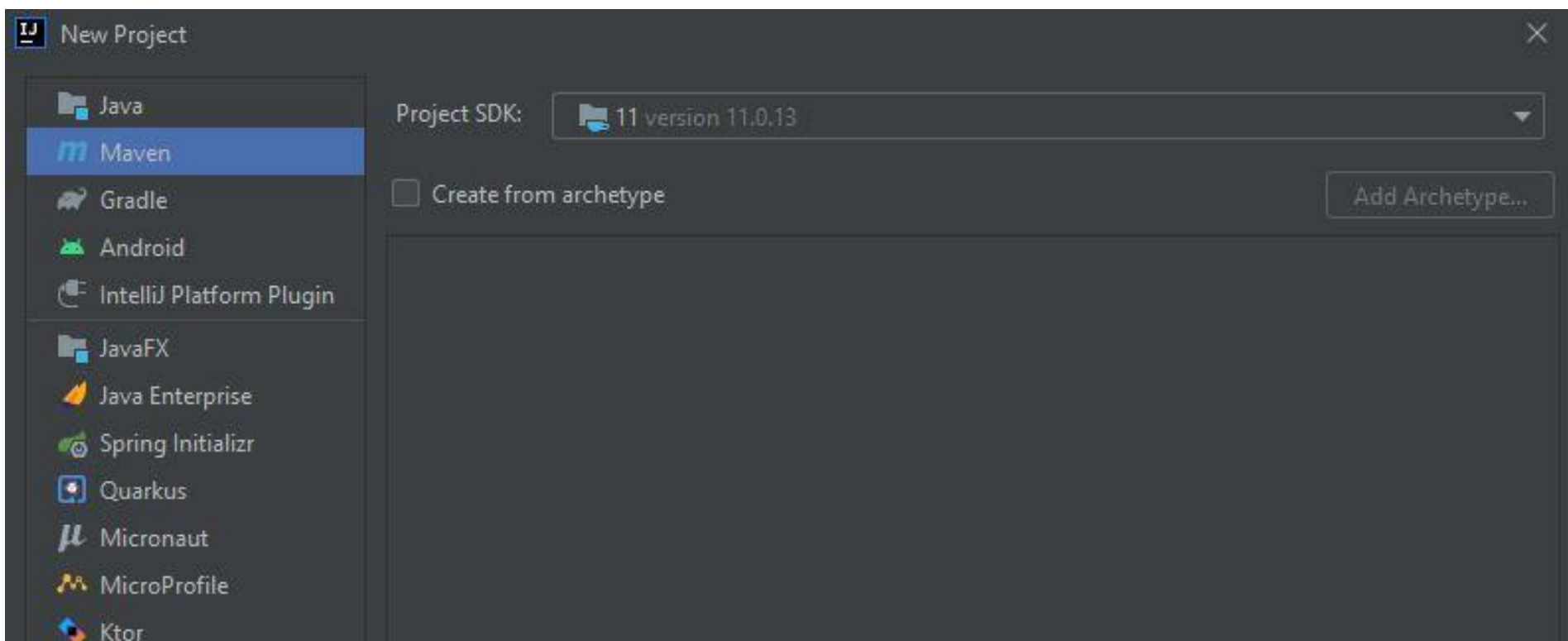
Installing from HDFS

```
DataFrame wordRDD = spark.Read().Text("hdfs://user/file");
```

Download from zipped file

```
DataFrame wordRDD = spark.Read().Text("*.gz");
```


SPARK USING JAVA



SPARK USING JAVA

mvnrepository.com/search?q=spark+core

Google ☆ Bookmarks PQ Arama Sonuçları - P... Academic conferen... Gelen Kutusu (11) -... International Journ... pysolar/solar.py at...

ITORY spark core Search

Found 36161 results

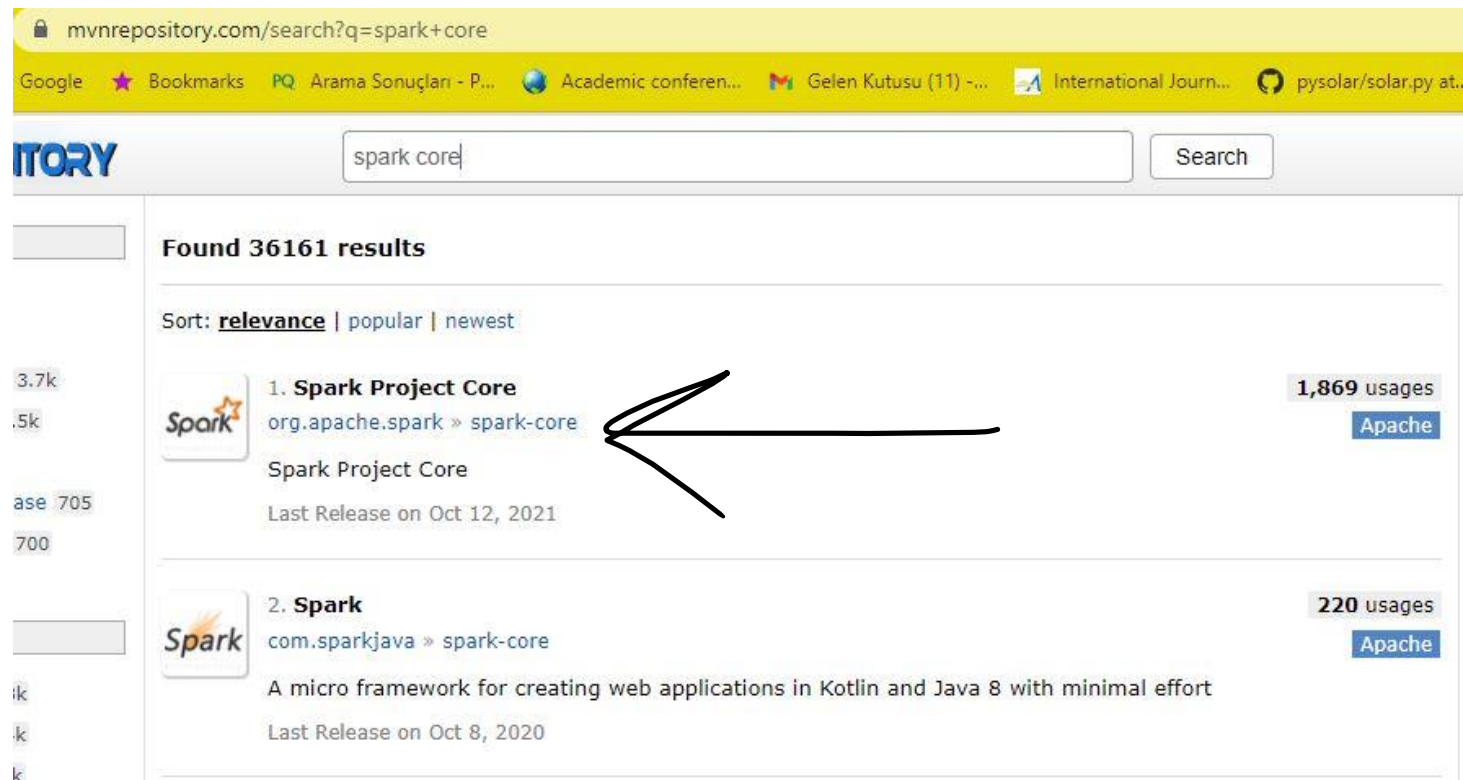
Sort: **relevance** | popular | newest

3.7k
.5k
ase 705
700

1. **Spark Project Core**
org.apache.spark » spark-core
Spark Project Core
Last Release on Oct 12, 2021
1,869 usages
Apache

2. **Spark**
com.sparkjava » spark-core
A micro framework for creating web applications in Kotlin and Java 8 with minimal effort
Last Release on Oct 8, 2020
220 usages
Apache

ik
k
k



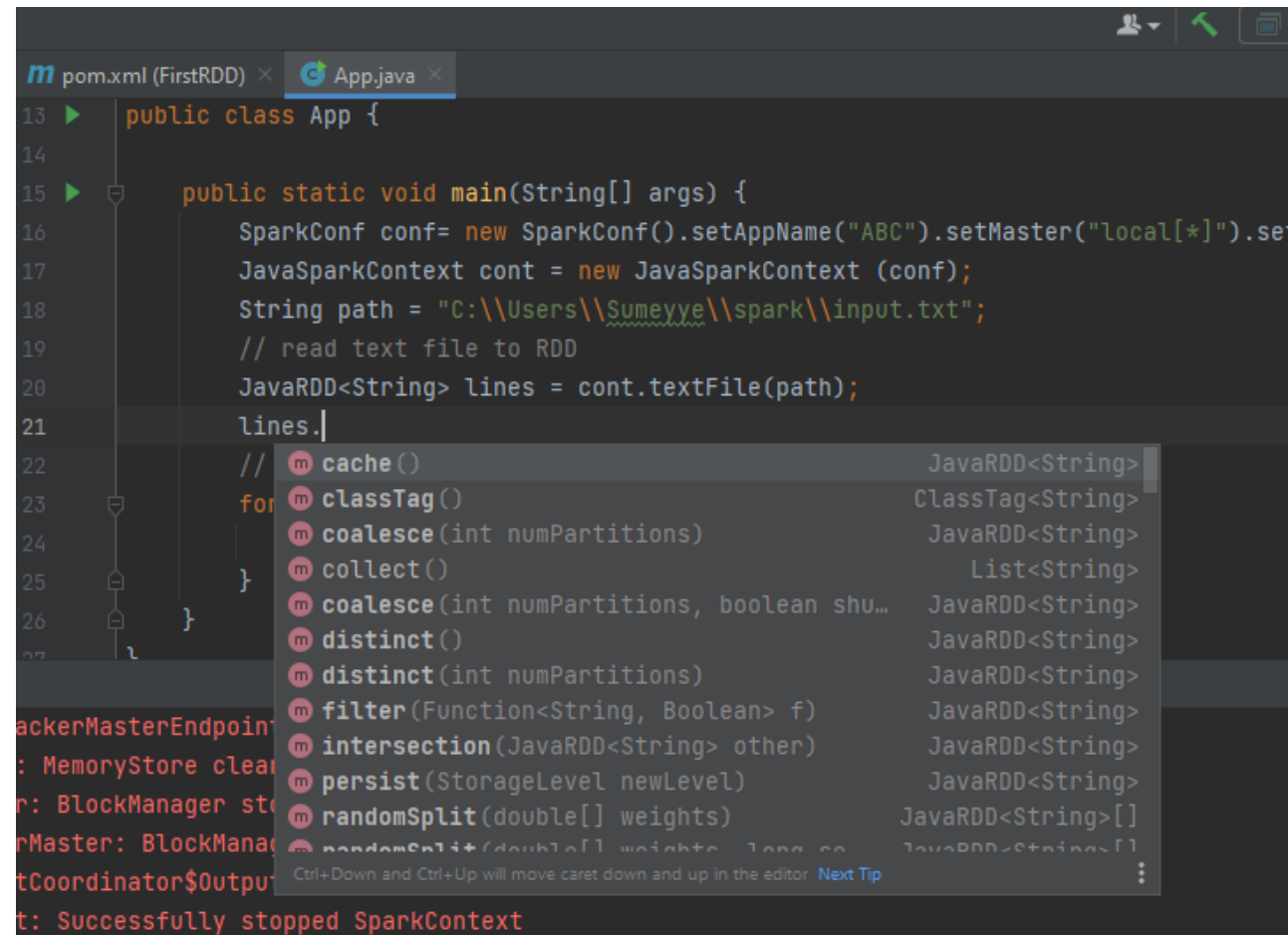
DEPENDENCE IN JAVA

```
<dependencies>  
  <!-- https://mvnrepository.com/artifact/org.apache.spark/spark-core -->  
  <dependency>  
    <groupId>org.apache.spark</groupId>  
    <artifactId>spark-core_2.12</artifactId>  
    <version>3.0.1</version>  
  </dependency>  
</dependencies>
```

CREATE FIRST RDD

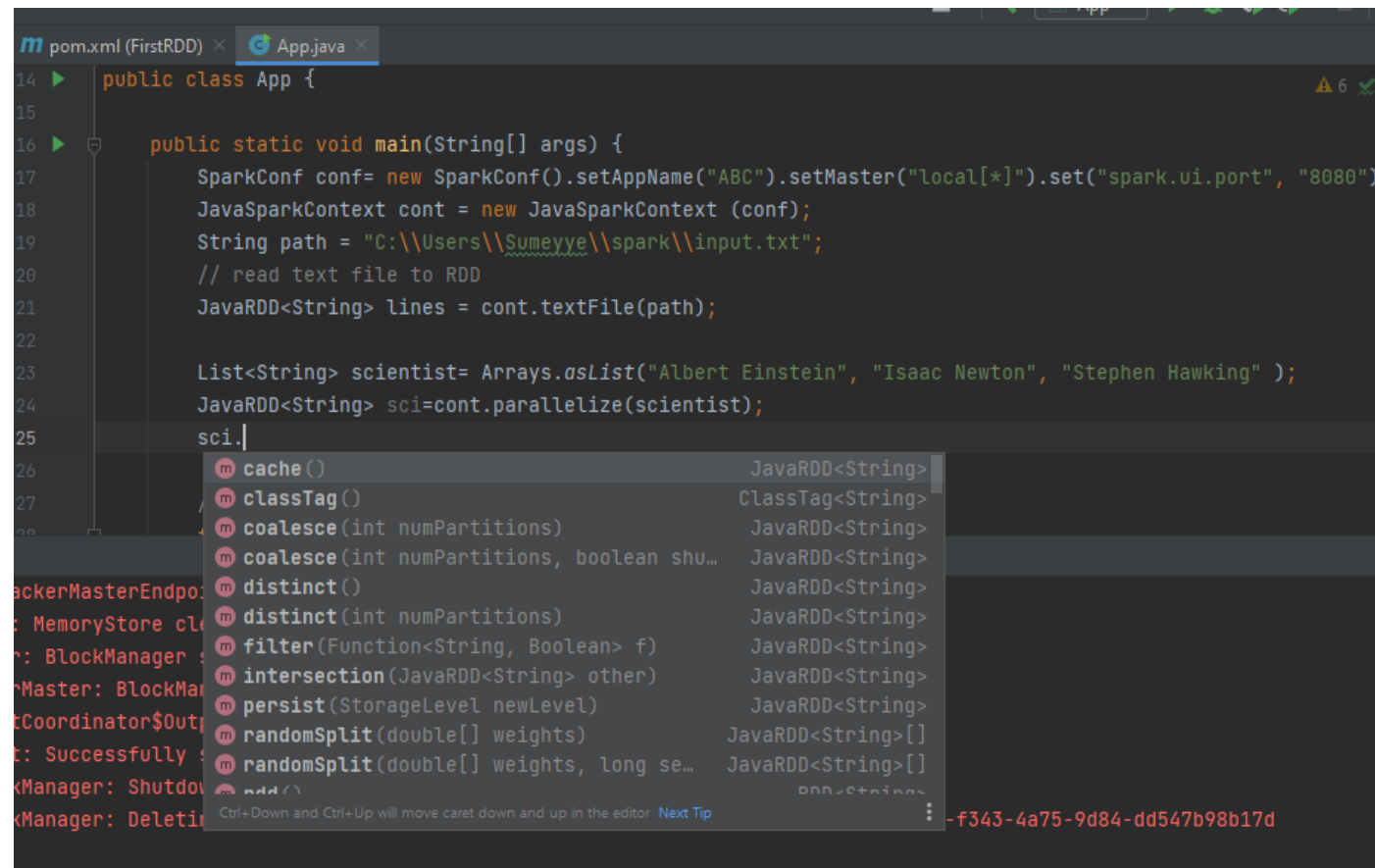
```
public static void main(String[] args) {  
    SparkConf conf= new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");;  
    JavaSparkContext sc = new JavaSparkContext (conf);  
    String path = "C:\\Users\\Sumeyye\\spark\\input.txt";  
    // read text file to RDD  
    JavaRDD<String> lines = sc.textFile(path);  
    // collect RDD for printing  
    for(String line:lines.collect()){  
        System.out.println(line);  
    }  
}
```

CREATE FIRST RDD



```
13 public class App {
14
15     public static void main(String[] args) {
16         SparkConf conf= new SparkConf().setAppName("ABC").setMaster("local[*]").set
17         JavaSparkContext cont = new JavaSparkContext (conf);
18         String path = "C:\\Users\\Sumeyye\\spark\\input.txt";
19         // read text file to RDD
20         JavaRDD<String> lines = cont.textFile(path);
21         lines.|
22         // m cache() JavaRDD<String>
23         for m classTag() ClassTag<String>
24         m coalesce(int numPartitions) JavaRDD<String>
25         m collect() List<String>
26         m coalesce(int numPartitions, boolean shu... JavaRDD<String>
27         m distinct() JavaRDD<String>
28         m distinct(int numPartitions) JavaRDD<String>
29         m filter(Function<String, Boolean> f) JavaRDD<String>
30         m intersection(JavaRDD<String> other) JavaRDD<String>
31         m persist(StorageLevel newLevel) JavaRDD<String>
32         m randomSplit(double[] weights) JavaRDD<String>[]
33         m randomSplit(double[] weights, long co JavaRDD<String>[]
34         Ctrl+Down and Ctrl+Up will move caret down and up in the editor Next Tip
35
36         t: Successfully stopped SparkContext
```

DATA LOAD TYPES

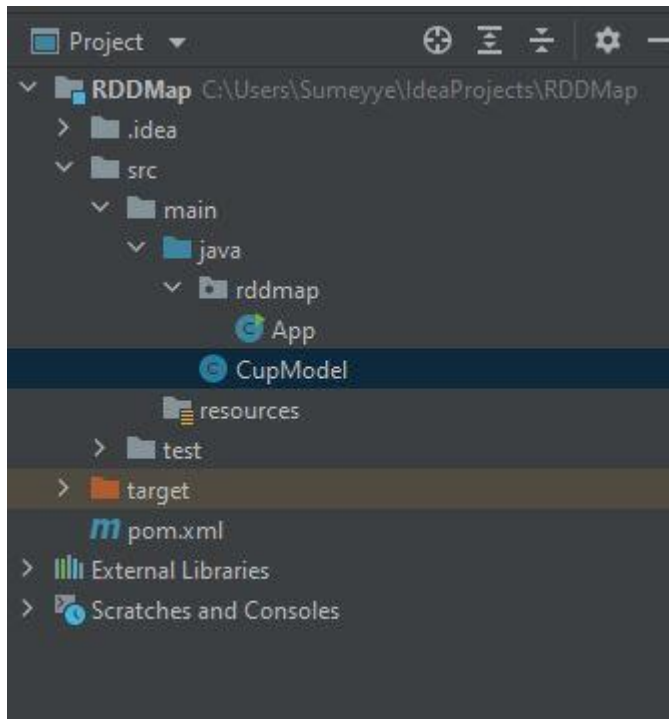


```
m pom.xml (FirstRDD) x App.java x
14 public class App {
15
16 public static void main(String[] args) {
17     SparkConf conf= new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");
18     JavaSparkContext cont = new JavaSparkContext (conf);
19     String path = "C:\\Users\\Sumeyye\\spark\\input.txt";
20     // read text file to RDD
21     JavaRDD<String> lines = cont.textFile(path);
22
23     List<String> scientist= Arrays.asList("Albert Einstein", "Isaac Newton", "Stephen Hawking" );
24     JavaRDD<String> sci=cont.parallelize(scientist);
25     sci.|
26     m cache() JavaRDD<String>
27     m classTag() ClassTag<String>
28     m coalesce(int numPartitions) JavaRDD<String>
29     m coalesce(int numPartitions, boolean shuffle) JavaRDD<String>
30     m distinct() JavaRDD<String>
31     m distinct(int numPartitions) JavaRDD<String>
32     m filter(Function<String, Boolean> f) JavaRDD<String>
33     m intersection(JavaRDD<String> other) JavaRDD<String>
34     m persist(StorageLevel newLevel) JavaRDD<String>
35     m randomSplit(double[] weights) JavaRDD<String>[]
36     m randomSplit(double[] weights, long seed) JavaRDD<String>[]
37     m rdd() RDD<String>
38     m ... RDD<String>
39     Ctrl+Down and Ctrl+Up will move caret down and up in the editor. Next Tip -f343-4a75-9d84-dd547b98b17d
```

RDD OPERATIONS

- Transformation: The result of operations on an rdd creates a new rdd. (map, filter)
- Actions: Calculation and saving operations on an rdd are done. (count, first)

CREATE JAVA OBJECT



```
import scala.Int;

public class CupModel {
    String year;
    String host;
    String first;
    String second;
    String third;
    String fourth;
    Int totalgoals;
    Int totalcountry;
    Int totalmatches;
    Int totalpartic;

    public CupModel(String year, String host, String first, String second, String third, String fourth,
        this.year = year;

    this.host = host;
    this.first = first;
    this.second = second;
    this.third = third;
    this.fourth = fourth;
    this.totalgoals = totalgoals;
    this.totalcountry = totalcountry;
    this.totalmatches = totalmatches;
    this.totalpartic = totalpartic;
}

public String getYear() {
    return year;
}

public void setYear(String year) {
```

```
public class App {  
    public static void main(String[] args) {  
        SparkConf conf= new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");  
        JavaSparkContext cont = new JavaSparkContext (conf);  
        String path = "C:\\bin\\WorldCup\\WorldCups.csv";  
        // read text file to RDD  
        JavaRDD<String> Raw_Data = cont.textFile(path);  
        System.out.println(Raw_Data.count());  
  
        Raw_Data.map(new Function<String, CupModel>() {  
            @Override  
            public CupModel call(String line) throws Exception {  
                return null;  
            }  
        });  
    }  
}
```

Raw_Data, String türündendir. Bu nedenle; new function ilk parametresi string olmalıdır. Bu String yapıyı object çevireceğiz. "CupModel" isminde object oluşturmuştuk. Fonksiyonu yazdıktan sonra implemente edilmelidir.

```

JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {
    @Override
    public CupModel call(String line) throws Exception {
        String[] split = line.split(",");
        var cupModel = new CupModel(split[0], split[1], split[2], split[3], split[4], split[5], Integer.par
        return cupModel;
    }
});

```

```

JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {
    @Override
    public CupModel call(String line) throws Exception {
        String[] split = line.split(regex: ",");
        var cupModel = new CupModel(split[0], split[1], split[2], split[3], split
        return cupModel;
    }
});

```

map.

- m map(Function<CupModel, R> f) JavaRDD<R>
 - m cache() JavaRDD<CupModel>
 - m classTag() ClassTag<CupModel>
 - m coalesce(int numPartitions) JavaRDD<CupModel>
 - m coalesce(int numPartitions, boolean s... JavaRDD<CupModel>
 - m distinct() JavaRDD<CupModel>
 - m distinct(int numPartitions) JavaRDD<CupModel>
 - m filter(Function<CupModel, Boolean> f) JavaRDD<CupModel>
 - m intersection(JavaRDD<CupModel> other) JavaRDD<CupModel>
 - m persist(StorageLevel newLevel) JavaRDD<CupModel>
 - m randomSplit(double[] weights) JavaRDD<CupModel>[]
 - m randomSplit(double[] weights, long JavaRDD<CupModel>[]
- Press Enter to insert, Tab to replace Next Tip

ava:27:107


```
JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {  
    @Override  
    public CupModel call(String line) throws Exception {  
        String[] split = line.split(regex: " ");  
        var cupModel = new CupModel(split[0], split[1], split[2], split[3], split[4], split[5], Integer.  
        return cupModel;  
    }  
});  
map.foreach(new VoidFunction<CupModel>() {  
    @Override  
    public void call(CupModel cupModel) throws Exception {  
        System.out.println(cupModel.getHost());  
    }  
});
```

```
public static void main(String[] args) {
    SparkConf conf= new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");
    JavaSparkContext cont = new JavaSparkContext (conf);
    String path = "C:\\bin\\WorldCup\\WorldCups.csv";
    // read text file to RDD
    JavaRDD<String> Raw_Data = cont.textFile(path);
    System.out.println(Raw_Data.count());

    JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {
        @Override
        public CupModel call(String line) throws Exception {
            String[] split = line.split(",");
            var cupModel = new CupModel(
                split[0],
                split[1],
                split[2],
                split[3],
                split[4],
                split[5]);
            return cupModel;
        }
    });
    map.foreach(new VoidFunction<CupModel>() {
        @Override
        public void call(CupModel cupModel) throws Exception {
            System.out.println(cupModel.getHost());
        }
    });
}
```

```
JavaRDD<CupModel> italy = map.filter(new Function<CupModel, Boolean>() {  
    @Override  
    public Boolean call(CupModel cupModel) throws Exception {  
        boolean italy = cupModel.getFirst().equals("Italy");  
        return italy;  
    }  
});  
italy.foreach(new VoidFunction<CupModel>() {  
    @Override  
    public void call(CupModel cupModel) throws Exception {  
        System.out.println(cupModel.getYear()+" "+cupModel.getFirst());  
    }  
});
```

```
21/12/14 09:44:25  
1982 Italy  
2006 Italy  
1934 Italy  
1938 Italy  
21/12/14 09:44:25
```

```
JavaRDD<String> flatmapRDD = Raw_Data.flatMap(new FlatMapFunction<String, String>() {  
    @Override  
    public Iterator<String> call(String s) throws Exception {  
        return Arrays.asList(s.split(regex: ",")).iterator();  
    }  
});  
  
flatmapRDD.foreach(new VoidFunction<String>() {  
    @Override  
    public void call(String s) throws Exception {  
        System.out.println(s);  
    }  
});
```

```
21/12/14 10:58:09 INFO HadoopRDD: Input s  
1982  
Spain  
Italy  
Germany FR  
Poland  
France  
146  
24  
52  
2.109.723  
1986  
Mexico
```

RDD OPERATIONS

- Transformation: The result of operations on an rdd creates a new rdd.
- (flatmap, mappartitions, distinct, sample, substract...)
- Actions: Calculation and saving operations on an rdd are done. (count, first)

```

public class Distinct {
    public static void main(String[] args) {
        SparkConf conf = new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");
        JavaSparkContext cont = new JavaSparkContext(conf);
        String path = "C:\\bin\\WorldCup\\WorldCups.csv";
        JavaRDD<String> Raw_Data = cont.textFile(path);

        JavaRDD<String> distinct = Raw_Data.distinct();
    }
}

```

```

public static void main(String[] args) {
    SparkConf conf = new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");
    JavaSparkContext cont = new JavaSparkContext(conf);
    String path = "C:\\bin\\WorldCup\\WorldCups.csv";
    JavaRDD<String> Raw_Data = cont.textFile(path);

    JavaRDD<String> sample = Raw_Data.sample(withReplacement: false, fraction: 0.5);
    sample.foreach(new VoidFunction<String>() {
        @Override
        public void call(String s) throws Exception {
            System.out.println(s);
        }
    });
}

```

```
public class Union {  
    public static void main(String[] args) {  
        SparkConf conf = new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");  
        JavaSparkContext cont = new JavaSparkContext(conf);  
        String first_path = "C:\\\\bin\\\\WorldCup\\\\WorldCups.csv";  
  
        JavaRDD<String> Raw_Data1 = cont.textFile(first_path);  
        System.out.println(Raw_Data1.count());  
  
        String second_path = "C:\\\\bin\\\\WorldCup\\\\WorldCupPlayers.csv";  
        JavaRDD<String> Raw_Data2 = cont.textFile(second_path);  
        System.out.println(Raw_Data2.count());  
  
        JavaRDD<String> unionsample = Raw_Data1.union(Raw_Data2);  
        System.out.println(unionsample.count());  
    }  
}
```


PAIRRDD

```
JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {  
    @Override  
    public CupModel call(String line) throws Exception {  
        String[] split = line.split(regex: " ");  
        var cupModel = new CupModel(  
            split[0],  
            split[1],  
            split[2],  
            split[3],  
            split[4],  
            split[5]);  
        return cupModel;  
    }  
});  
  
JavaPairRDD<String, String> JavaPairRDD = map.mapToPair(new PairFunction<CupModel, String, String>() {  
    @Override  
    public Tuple2<String, String> call(CupModel cupModel) throws Exception {  
        return new Tuple2<String, String>(cupModel.getFirst(), cupModel.getSecond());  
    }  
});
```

```
JavaPairRDD.foreach(new VoidFunction<Tuple2<String, String>>() {  
    @Override  
    public void call(Tuple2<String, String> line) throws Exception {  
        System.out.println(line._1 + " " + line._2);  
    }  
});
```

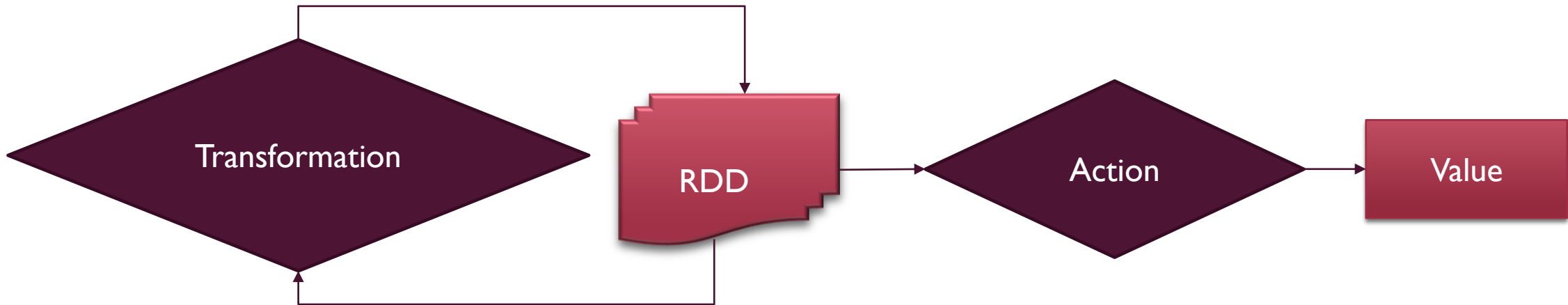

PAIRRDD

```
org.apache.spark.api.java.JavaPairRDD<String, Iterable<String>> result = JavaPairRDD.groupByKey();

result.foreach(new VoidFunction<Tuple2<String, Iterable<String>>>() {
    @Override
    public void call(Tuple2<String, Iterable<String>> line) throws Exception {
        System.out.println(line._1+ " " + line._2);
    }
});
```


LAZY EVALUATION

- Transformation işlemlerinde herhangi bir işlem yapılmaz. Spark, action metodunu görünce işlem başlatır. Spark, action yöntemini gördüğünde değere dönüştürür.



LAZY EVALUATION

```
JavaRDD<CupModel> map = Raw_Data.map(new Function<String, CupModel>() {  
    @Override  
    public CupModel call(String line) throws Exception {  
        System.out.println(line);  
        String[] split = line.split(regex: ",");  
        var cupModel = new CupModel(  
            split[0],  
            split[1],  
            split[2],  
            split[3],  
            split[4],  
            split[5]);  
        return cupModel;  
    }  
});  
  
System.out.println(map.count());
```

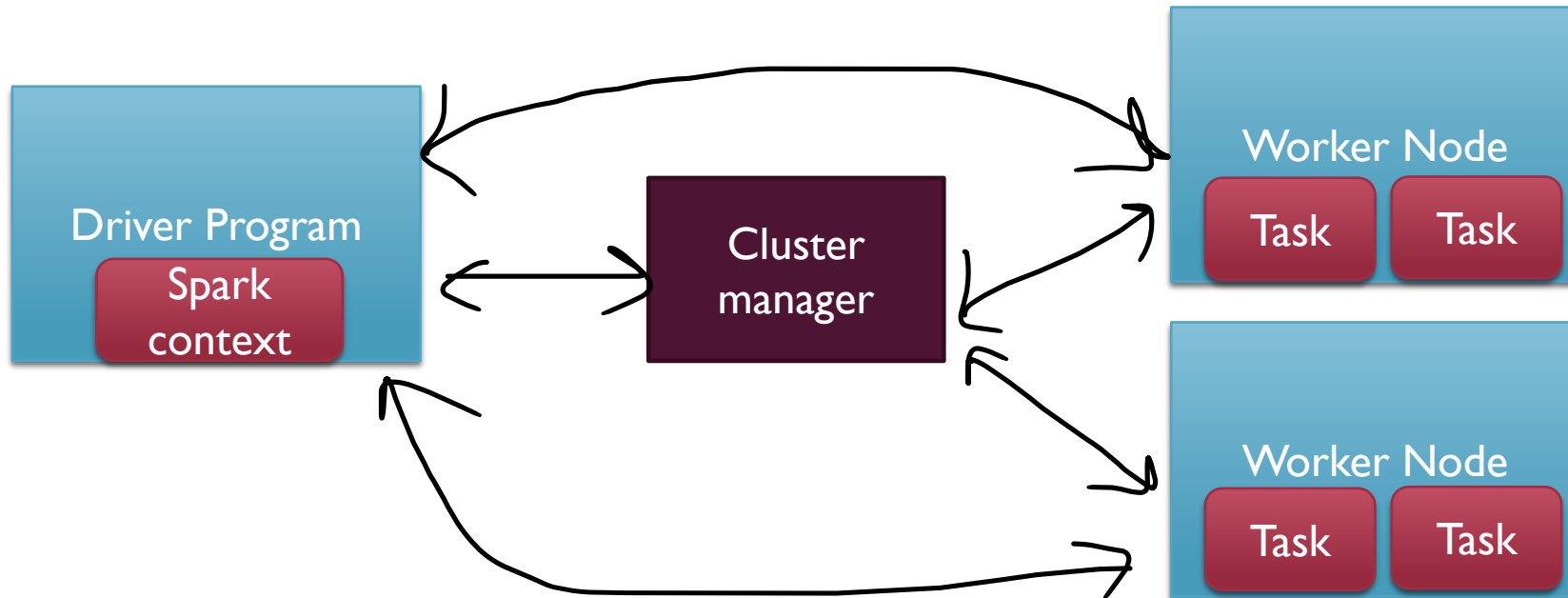


ACTION METHOD

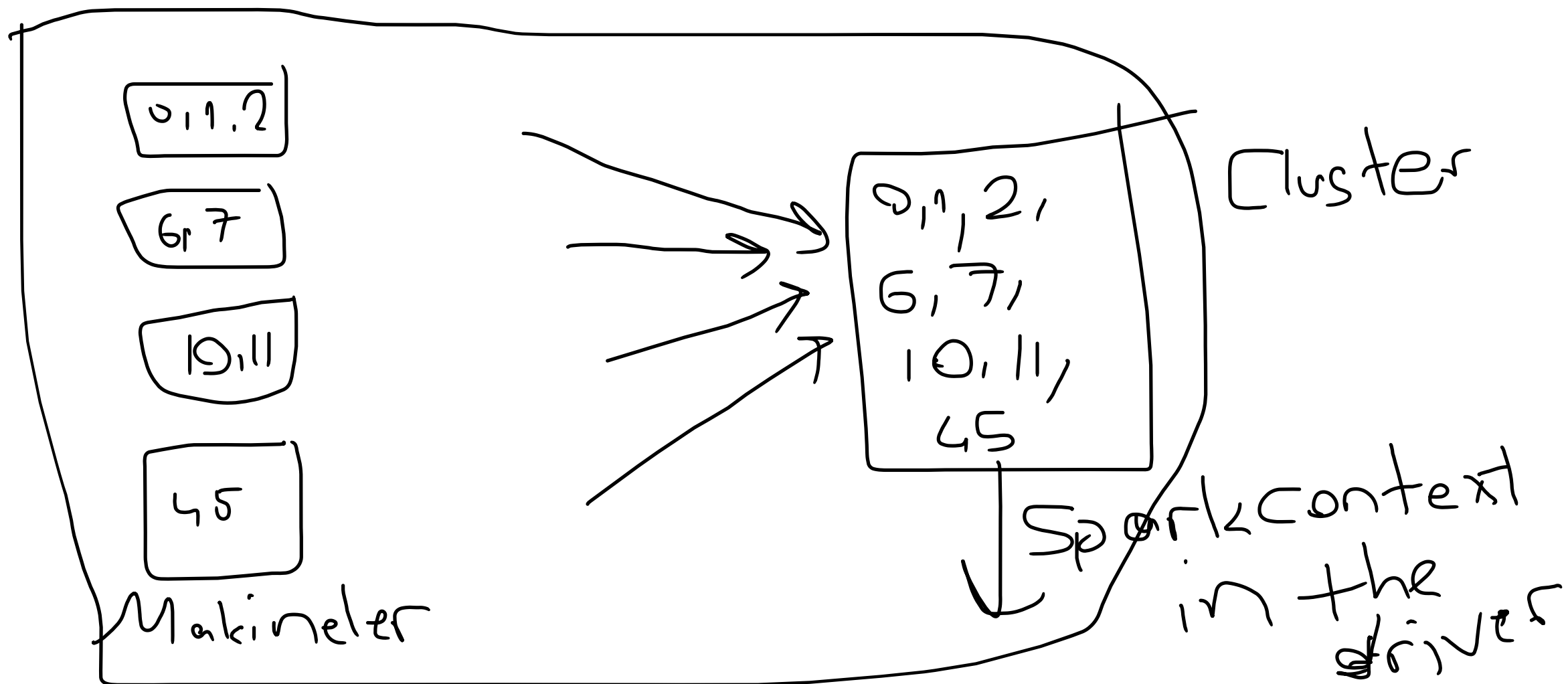
- Count
- First
- Collect
- Foreach
- Filter
- Take vb.

COLLECT

- It sends the data on the machines in the cluster to the master machine.



COLLECT



COLLECT

```
public class Collect {  
    public static void main(String[] args) {  
        SparkConf conf = new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");  
        JavaSparkContext cont = new JavaSparkContext(conf);  
        String path = "C:\\\\bin\\\\WorldCup\\\\WorldCups.csv";  
        // read text file to RDD  
        JavaRDD<String> Raw_Data = cont.textFile(path);  
  
        for (String s : Raw_Data.collect()) {  
            System.out.println(s);  
        }  
    }  
}
```

TAKE

```
public static void main(String[] args) {  
    SparkConf conf = new SparkConf().setAppName("ABC").setMaster("local[*]").set("spark.ui.port", "8080");  
    JavaSparkContext cont = new JavaSparkContext(conf);  
    String path = "C:\\\\bin\\\\WorldCup\\\\WorldCups.csv";  
    // read text file to RDD  
    JavaRDD<String> Raw_Data = cont.textFile(path);  
    for (String s : Raw_Data.collect()) {  
        System.out.println(s);  
    }  
    List<String> take = Raw_Data.take(4);  
}
```


TAKESAMPLE-SAVEASTEXTFILE

```
List<String> take = Raw_Data.take( num: 4);  
List<String> takesample = Raw_Data.takeSample( withReplacement: false, num: 5);
```

```
Raw_Data.saveAsTextFile( path: "C:\\bin\\Raw_Data");
```