Session 25:
BIG DATA ECOSYSTEM
NTEGRATION
Assignment 1
-Prachi Mohite

## Task 1
## As discussed in class integrate Spark Hive

To demonstrate this we have a program which will list the databases in HIVE

Steps to be followed
1. Copy the hive-site.xml file from $HIVE_HOME/conf to $SPARK_HOME/conf



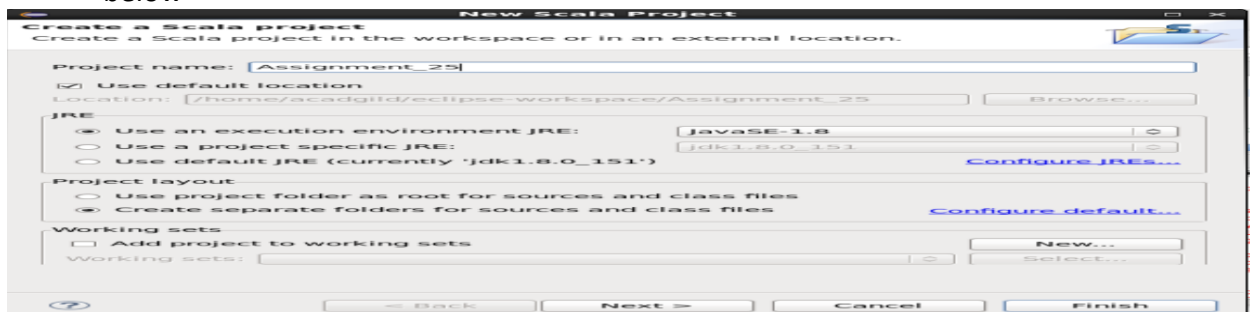2. Add the following properties to hive-site.xml on spark side :

   &lt;property&gt;

   &lt;name&gt;hive.metastore.uris&lt;/name&gt;

   &lt;value&gt;thrift://localhost:9083&lt;/value&gt;

   &lt;description&gt;URI Client to connect to metastore service&lt;/description&gt;

   &lt;/property&gt;

```
<property>
        <name>hive.metastore.uris</name>
        <value>thrift://localhost:9083</value>
        <description>URI for client to connect to metastore service</description>
</property>

</configuration>
```
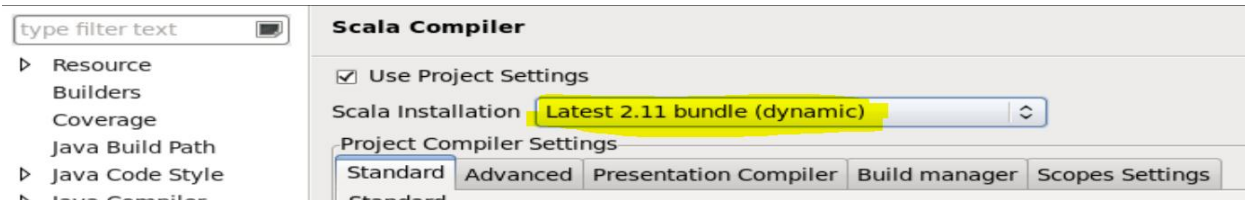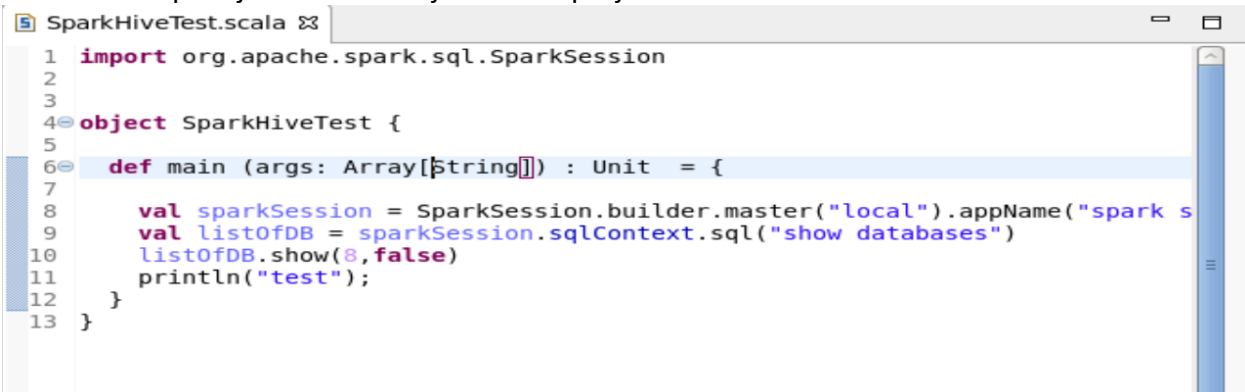
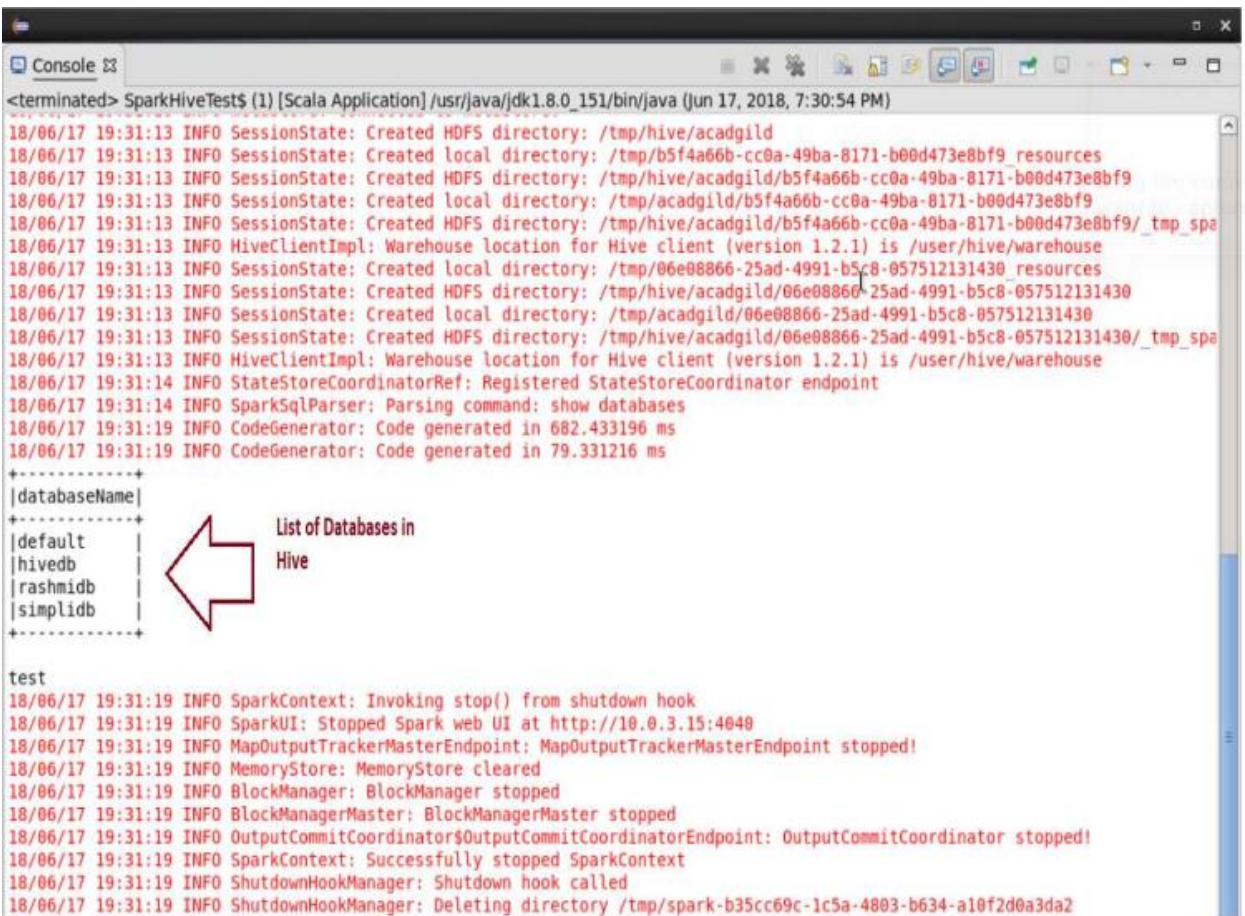3. Downloaded the code from here and added to a project created as Assignment_25 as below



4. After adding the file to project change the scala compiler to Scala 2.11 as below

Scala Compiler

☑ Use Project Settings

Scala Installation  Latest 2.11 bundle (dynamic)  ⬍

Project Compiler Settings

| Standard | Advanced | Presentation Compiler | Build manager | Scopes Settings |

5. Add Spark jars and hbase jars to the project

```scala
1  import org.apache.spark.sql.SparkSession
2
3
4  object SparkHiveTest {
5
6    def main (args: Array[String]) : Unit  = {
7
8      val sparkSession = SparkSession.builder.master("local").appName("spark s
9      val listOfDB = sparkSession.sqlContext.sql("show databases")
10     listOfDB.show(8,false)
11     println("test");
12   }
13 }
```

6. Execute the code and verify the list of databases

```
Console ⬚

<terminated> SparkHiveTest$ (1) [Scala Application] /usr/java/jdk1.8.0_151/bin/java (Jun 17, 2018, 7:30:54 PM)
18/06/17 19:31:13 INFO SessionState: Created HDFS directory: /tmp/hive/acadgild
18/06/17 19:31:13 INFO SessionState: Created local directory: /tmp/b5f4a66b-cc0a-49ba-8171-b00d473e8bf9_resources
18/06/17 19:31:13 INFO SessionState: Created HDFS directory: /tmp/hive/acadgild/b5f4a66b-cc0a-49ba-8171-b00d473e8bf9
18/06/17 19:31:13 INFO SessionState: Created local directory: /tmp/acadgild/b5f4a66b-cc0a-49ba-8171-b00d473e8bf9
18/06/17 19:31:13 INFO SessionState: Created HDFS directory: /tmp/hive/acadgild/b5f4a66b-cc0a-49ba-8171-b00d473e8bf9/_tmp_spa
18/06/17 19:31:13 INFO HiveClientImpl: Warehouse location for Hive client (version 1.2.1) is /user/hive/warehouse
18/06/17 19:31:13 INFO SessionState: Created local directory: /tmp/06e08866-25ad-4991-b5c8-057512131430_resources
18/06/17 19:31:13 INFO SessionState: Created HDFS directory: /tmp/hive/acadgild/06e08866-25ad-4991-b5c8-057512131430
18/06/17 19:31:13 INFO SessionState: Created local directory: /tmp/acadgild/06e08866-25ad-4991-b5c8-057512131430
18/06/17 19:31:13 INFO SessionState: Created HDFS directory: /tmp/hive/acadgild/06e08866-25ad-4991-b5c8-057512131430/_tmp_spa
18/06/17 19:31:13 INFO HiveClientImpl: Warehouse location for Hive client (version 1.2.1) is /user/hive/warehouse
18/06/17 19:31:14 INFO StateStoreCoordinatorRef: Registered StateStoreCoordinator endpoint
18/06/17 19:31:14 INFO SparkSqlParser: Parsing command: show databases
18/06/17 19:31:19 INFO CodeGenerator: Code generated in 682.433196 ms
18/06/17 19:31:19 INFO CodeGenerator: Code generated in 79.331216 ms
+------------+
|databaseName|
+------------+
|default     |
|hivedb      |
|rashmidb    |
|simplidb    |
+------------+

test
18/06/17 19:31:19 INFO SparkContext: Invoking stop() from shutdown hook
18/06/17 19:31:19 INFO SparkUI: Stopped Spark web UI at http://10.0.3.15:4040
18/06/17 19:31:19 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/06/17 19:31:19 INFO MemoryStore: MemoryStore cleared
18/06/17 19:31:19 INFO BlockManager: BlockManager stopped
18/06/17 19:31:19 INFO BlockManagerMaster: BlockManagerMaster stopped
18/06/17 19:31:19 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/06/17 19:31:19 INFO SparkContext: Successfully stopped SparkContext
18/06/17 19:31:19 INFO ShutdownHookManager: Shutdown hook called
18/06/17 19:31:19 INFO ShutdownHookManager: Deleting directory /tmp/spark-b35cc69c-1c5a-4803-b634-a10f2d0a3da2
```

List of Databases in Hive

## Task 2
## As discussed in class integrate Spark HBase

As we have already done some pre-requisite steps. Now download the code and add to the project. This code will integrate with HBase through spark and create a table and insert contents in the same.
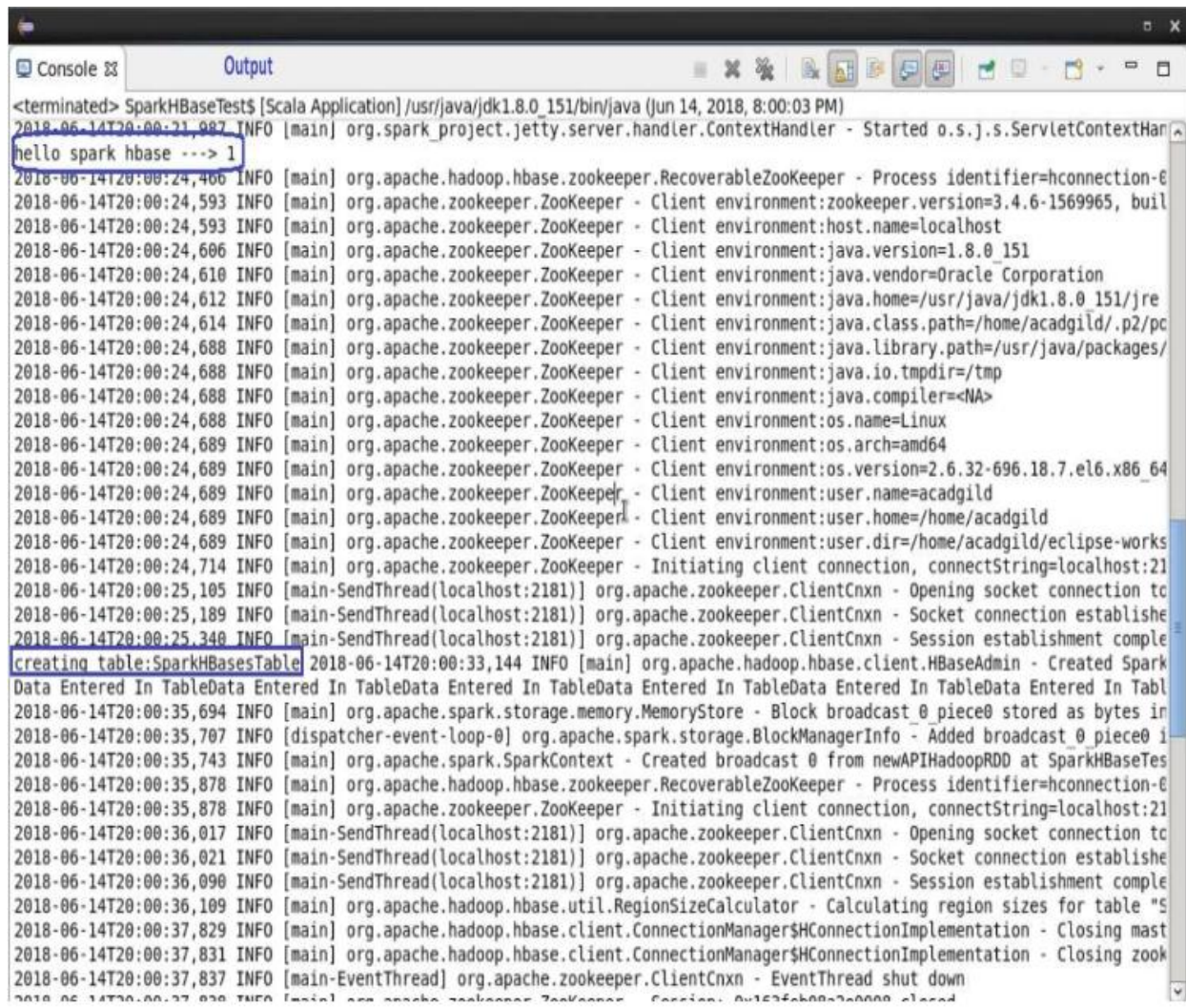
Add the required jar files for the HBase.

Make sure HBase is started while executing the code.

Added Code as below

```scala
S *SparkHBaseTest.scala ⋈

15    def main(args: Array[String]) {
16        // Create a SparkContext using every core of the local machine, named RatingsCounter
17        val sc = new SparkContext("local[*]", "SparkHBaseTest")
18
19        println("hello spark hbase ---> 1")
20
21        val conf = HBaseConfiguration.create()
22        val tablename = "SparkHBasesTable"
23        conf.set(TableInputFormat.INPUT_TABLE,tablename)
24        val admin = new HBaseAdmin(conf)
25        if(!admin.isTableAvailable(tablename)){
26            print("creating table:"+tablename+"\t")
27            val tableDescription = new HTableDescriptor(tablename)
28            tableDescription.addFamily(new HColumnDescriptor("cf".getBytes()));
29            admin.createTable(tableDescription);
30        } else {
31            print("table already exists")
32        }
33
34        val table = new HTable(conf,tablename);
35        for(x <- 1 to 10){
36            var p = new Put(new String("row" + x).getBytes());
37            p.add("cf".getBytes(),"column1".getBytes(),new String("value" + x).getBytes());
38            table.put(p);
39            print("Data Entered In Table")
40        }
41        val hBaseRDD = sc.newAPIHadoopRDD(conf, classOf[TableInputFormat],
42            classOf[ImmutableBytesWritable],classOf[Result])
43        print("RecordCount->>"+hBaseRDD.count())
44        sc.stop()
```

**Execution of Code**

Hre before execution of the Code, we cannot see the table is created. However after execution of the code we could see table is created and records are inserted into it as below.

```
hbase(main):005:0> list
TABLE                               Before executing Spark HBase integration program
0 row(s) in 0.0340 seconds

=> []
hbase(main):006:0> list
TABLE                               After executing Spark HBase integration program
SparkHBasesTable
1 row(s) in 0.0370 seconds
                                    Contents of the HBase table created i.e. SparkHBaseTable
=> ["SparkHBasesTable"]
hbase(main):007:0> scan 'SparkHBasesTable'
ROW                      COLUMN+CELL
 row1                    column=cf:column1, timestamp=1528992429713, value=value1
 row10                   column=cf:column1, timestamp=1528992429816, value=value10
 row2                    column=cf:column1, timestamp=1528992429746, value=value2
 row3                    column=cf:column1, timestamp=1528992429757, value=value3
 row4                    column=cf:column1, timestamp=1528992429764, value=value4
 row5                    column=cf:column1, timestamp=1528992429771, value=value5
 row6                    column=cf:column1, timestamp=1528992429777, value=value6
 row7                    column=cf:column1, timestamp=1528992429785, value=value7
 row8                    column=cf:column1, timestamp=1528992429800, value=value8
 row9                    column=cf:column1, timestamp=1528992429807, value=value9
10 row(s) in 0.1380 seconds

hbase(main):008:0>
```

**Task 3**
**As discussed in class integrate Spark HBase**

Pre-requisite

Start the zookeeper server in Kafka by navigating into $KAFKA_HOME with the command given below:

```
./bin/zookeeper-server-start.sh config/zookeeper.properties
```

Keep the terminal running, open one new terminal, and start the Kafka broker using the following command:

```
./bin/kafka-server-start.sh config/server.properties
```

After starting, leave both the terminals running, open a new terminal, and create a Kafka topic with the following command:

```
./bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-
factor 1 --partitions 1 --topic acdgild-topic
```

```
[acadgild@localhost ~]$ cd $KAFKA_HOME
[acadgild@localhost kafka_2.12-0.10.1.1]$ ./bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-factor 1 --partitions 1 --topic acdgild-topic
Created topic "acdgild-topic".
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost kafka_2.12-0.10.1.1]$ ./bin/kafka-topics.sh --list --zookeeper localhost:2181
ItemTopic
KeyLessTopic
KeyedTopic
UserTopic
__consumer_offsets
acdgild-topic
[acadgild@localhost kafka_2.12-0.10.1.1]$
```

Program which runs the word count program by reading the contents from kafka and run in spark.

**Code**

```scala
 1  import org.apache.spark._
 2  import org.apache.spark.streaming.StreamingContext
 3  import org.apache.spark.streaming.Seconds
 4  import org.apache.spark.streaming.kafka.KafkaUtils
 5  object WordCount {
 6    def main( args:Array[String] ){
 7      val conf = new SparkConf().setMaster("local[*]").setAppName
 8      val ssc = new StreamingContext(conf, Seconds(10))
 9      val kafkaStream = KafkaUtils.createStream(ssc, "localhost:2
10  //need to change the topic name and the port number accordingly
11      val words = kafkaStream.flatMap(x =>  x._2.split(" "))
12      val wordCounts = words.map(x => (x, 1)).reduceByKey(_ + _)
13      kafkaStream.print()  //prints the stream of data received
14      wordCounts.print()   //prints the wordcount result of the s
15      ssc.start()
16      ssc.awaitTermination()
17    }
18  }
```

**Output**

```
[acadgild@localhost ~]$ cd $KAFKA_HOME
[acadgild@localhost kafka_2.12-0.10.1.1]$ ./bin/kafka-console-producer.sh --broker-list localhost:9092 --topic acadgild-topic
Hello,
This is BDH session. This is a wonderful Session.
This is a great session
great session wonderful session

Hello,
This is BDH session. This is a wonderful Session.
This is a great session
great session wonderful session
```

• Data inputted by user

```
Console ☒
JavaDirectKafkaWordCount [Java Application] /usr/java/jdk1.8.0_151/bin/java (Jun 14, 2018, 5:17:06 PM)
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/spark/spark-2.2.1-bin-hadoop
SLF4J: Found binding in [jar:file:/home/acadgild/install/kafka/kafka_2.12-0.10.1.1/li
SLF4J: Found binding in [jar:file:/home/acadgild/Downloads/jar_files(7)/slf4j-log4j12
SLF4J: Found binding in [jar:file:/home/acadgild/Downloads/jar_files(8)/slf4j-log4j12
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
-------------------------------------------
Time: 1528976850000 ms
-------------------------------------------
(Session.,1)
(is,3)
(session.,1)
(BDH,1)
(wonderful,2)
(session,3)
(This,3)
(Hello,,1)
(a,2)
(great,2)
```

Word Count is permormed on "acadgild-topic"