```
import org.apache.spark.SparkConf
 import org.apache.spark.SparkContext
 import org.apache.spark.sql.SparkSession
 import org.apache.spark.sql.types._
 import org.apache.spark.sql.Row
 import org.apache.spark.sql.functions.
 import org.apache.spark._
 import scala.io.Source
 import com.mysql.jdbc.Driver
⊖ object Task1 {
  def main(args:Array[String]):Unit={
       val conf = new SparkConf().setAppName("zipcode").setMaster("local[*]")
           val sc = new SparkContext(conf)
           sc.setLogLevel("Error")
           val spark = SparkSession.builder().getOrCreate()
           import spark.implicits.
           //Get Buffer response
           val http_buffer = Source.fromURL("https://randomuser.me/api/0.8/?results=10")
           //Create String from buffer
           val http_string = http_buffer.mkString
           //Create rdd[String] from String
           val rdd_str = sc.parallelize(List(http_string))
           //Convert rdd to data frame
           val df http = spark.read.json(rdd str)
           //df http.show()
```

```
//Convert rdd to data frame
val df http = spark.read.json(rdd str)
//df http.show()
df http.printSchema()
val df_http_all = df_http.withColumn("results", explode(col("results")))
.select("results.user.gender", "results.user.name.title", "results.user.name.first",
    "results.user.name.last", "results.user.location.street",
    "results.user.location.city", "results.user.location.state",
    "results.user.location.zip", "results.user.email",
    "results.user.username", "results.user.password",
    """
     "results.user.salt", "results.user.md5",
     "results.user.sha1", "results.user.sha256", "results.user.registered", "results.user.dob" results.user.phone", "results.user.cell",
     "results.user.picture.large",
"results.user.picture.medium", "results.user.picture.thumbnail",
"nationality", "seed", "version")
//write to hive
df_http_all.write.mode("overwrite").format("orc").saveAsTable("zeyodb.userdata")
println("data written to hive successfully")
val df_hive_read = spark.sql("select * from zeyodb.userdata")
df_hive_read.show()
println("data read from hive successfully")
df http all.write.format("jdbc").option("url","jdbc:mysql://localhost/zeyodb")
.mode("overwrite")
.option("driver", "com.mysql.jdbc.Driver")
.option("dbtable", "userdata")
.option("user","dbuser")
```

```
//write to hive
df_http_all.write.mode("overwrite").format("orc").saveAsTable("zeyodb.userdata")

println("data written to hive successfully")

val df_hive_read = spark.sql("select * from zeyodb.userdata")
df_hive_read.show()

println("data read from hive successfully")

df_http_all.write.format("jdbc").option("url","jdbc:mysql://localhost/zeyodb")
.mode("overwrite")
.option("driver", "com.mysql.jdbc.Driver")
.option("dbtable", "userdata")
.option("user", "dbuser")
.option("password", "cloudera").save()
println("Data Written to RDBMS")

val df_read_rdbms = spark.read.format("jdbc").option("url", "jdbc:mysql://localhost/zeyodb
.option("driver", "com.mysql.jdbc.Driver")
.option("driver", "com.mysql.jdbc.Driver")
.option("dbtable", "userdata")
.option("dbtable", "userdata")
.option("password", "cloudera").load()

df_read_rdbms.write.mode("overwrite").format("parquet").saveAsTable("zeyodb.parquet")
println("Data written to hive")
}
```

```
//Create rdd[String] from String
                      val rdd_str = sc.parallelize(List(http_string))
                      //Convert rdd to data frame
                      val df_http = spark.read.json(rdd_str)
                      //df_http.show()
                      df_http.printSchema()
                      val df http all = df http.withColumn("results", explode(col("results")))
                       .select("results.user.location.zip")
                      //df_http_all.show()
                      //Overwrite for the first time and append other times
                      var write mode = ""
                      if(i == 1)
                           write_mode = "overwrite"
                      else
                          write_mode = "append"
                      df http_all.write.format("jdbc").option("url","jdbc:mysql://localhost/zeyodb")
                      .mode(write mode)
                      .option("driver", "com.mysql.jdbc.Driver")
.option("dbtable","zipcode")
                      .option("user","dbuser")
                      .option("password","cloudera").save()
println("Written - " + i)
                                                                                                                     🦹 Problems 🏿 🗐 Tasks 📮 Console 🛭
<terminated> Task2$ [Scala Application] C:\Program Files\Java\jdk1.8.0_172\bin\javaw.exe (12-Jun-2020, 8:49:38 AM)
                  -- sha1: string (nullable = true)
                  -- sha256: string (nullable = true)
                  |-- username: string (nullable = true)
 |-- seed: string (nullable = true)
 |-- version: string (nullable = true)
Written - 10
 Administrator: Command Prompt - mysql -u dbuser -p
                                                                                                                                             90097
87071
93902
  34441
79059
52816
  75107
48591
  24021
  17448
51854
  86805
43392
62958
  19780
32671
  59670
  65765
27220
  79272
34826
63728
  58609
46654
  66121
 100 rows in set (0.00 sec)
```

```
import org.apache.spark.SparkConf
import org.apache.spark.SparkContext
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.types._
import org.apache.spark.sql.Row
import org.apache.spark.sql.functions._
import org.apache.spark._
import scala.io.Source
import com.mysql.jdbc.Driver
object Task3 {
  def main(args:Array[String]):Unit={
     val conf = new SparkConf().setAppName("zipcode").setMaster("local[*]")
         val sc = new SparkContext(conf)
         sc.setLogLevel("Error")
         val spark = SparkSession.builder().enableHiveSupport().getOrCreate()
         import spark.implicits._
         val repeat count = 10
         for(i <- 1 to repeat count)</pre>
                //Get Buffer response
                val http buffer = Source.fromURL("https://randomuser.me/api/0.8/?results=10")
                //Create String from buffer
                val http_string = http_buffer.mkString
                //Create rdd[String] from String
                val rdd_str = sc.parallelize(List(http_string))
                //Convert rdd to data frame
                val df_http = spark.read.json(rdd_str)
                //df_http.show()
```

```
val rdd_str = sc.parallelize(List(http_string))
                  //Convert rdd to data frame
                  val df_http = spark.read.json(rdd_str)
                  //df http.show()
                  df http.printSchema()
                  val df http all = df http.withColumn("results", explode(col("results")))
                  .select("results.user.location.zip")
                  //df http all.show()
                  //Overwrite for the first time and append other times
                  var write_mode = ""
                  if(i == 1)
                      write_mode = "overwrite"
                      write_mode = "append"
                  //Write to RDBMS
                  df http all.write.format("jdbc").option("url","jdbc:mysql://localhost/zeyodb")
                  .mode(write mode)
                  .option("driver", "com.mysql.jdbc.Driver")
.option("dbtable", "zipcode")
                  .option("user","dbuser")
                  .option("password","cloudera").save()
                  println("Written - " + i)
                  println("data written to rdbms")
                  //Write to Hive
                  df_http_all.write.format("hive").mode(write_mode).saveAsTable("zipcode")
                  println("data written to hive")
⊕ }
```

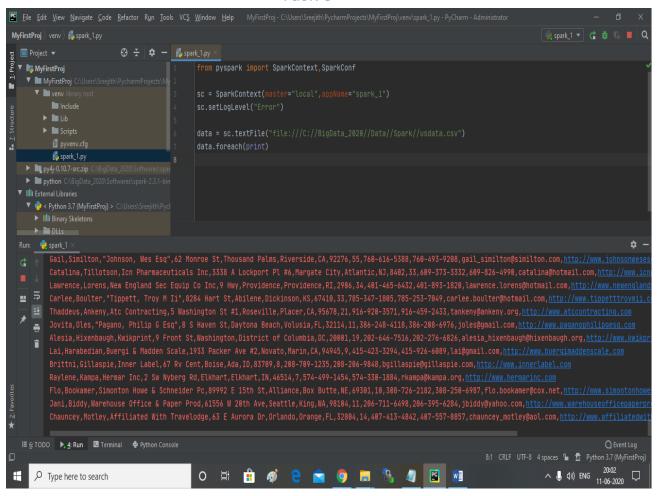
Command to submit job:

spark-submit --class Task.Task --master local[*] --jars "/home/cloudera/ext_jars/*" --conf spark.sql.catalogImplementation=hive /home/cloudera/program_jars/Task-0.0.1-SNAPSHOT.jar

Task 4

Scheduling the job using cron tab:

```
/20 * * * * spark-submit --class Task.Task --master local[*] --jars "/home/cloudera/ext_jars/*" --conf spark.sql.catalogImplementation=hive /home/cloudera/program jars/Task-0.0.1-SNAPSHOT.jar
```



```
⊖ object Task6 {
  def main(args:Array[String]):Unit={
          val conf = new SparkConf().setAppName("DSL").setMaster("local[*]")
              val sc = new SparkContext(conf)
              sc.setLogLevel("Error")
              val spark = SparkSession.builder().getOrCreate()
              import spark.implicits._
              val batters df = spark.read.option("multiline",true).format("json").load("file:///C://BigData 2020//Data//Spark/batters.
              //batters df.show()
              batters df.printSchema()
              val batters_select = batters_df.withColumn("batter", explode(col("batters.batter")))
                                             .withColumn("topping", explode(col("topping")))
                                             .select(col("id"),col("name"),col("ppu"),col("type"),col("batter.id").alias("batter_id"),
                                             col("batter.type").alias("batter_type"),col("topping.id").alias("topping_id"),
                                             col("topping.type").alias("topping type"))
              batters_select.show()
              batters_select.printSchema()
    }
    <
                                                                                                         🦹 Problems 🏿 🖪 Tasks 📮 Console 🛭
<terminated> Task6$ [Scala Application] C:\Program Files\Java\jdk1.8.0_172\bin\javaw.exe (12-Jun-2020, 8:04:31 AM)
  id|name| ppu| type|batter_id|batter_type|topping_id|
                                                              topping_type|
|0001|Cake|0.55|donut|
                          1001
                                   Regular
                                                 5001
                                                                      None
|0001|Cake|0.55|donut|
                          1001
                                   Regular
                                                 5002
                                                                    Glazed
|0001|Cake|0.55|donut|
                          1001
                                   Regular
                                                 5005
                                                                     Sugar
                                                            Powdered Sugar
|0001|Cake|0.55|donut|
                          1001
                                   Regular
                                                 5007
<
```