# Load data from AWS RDS to Hadoop

**Command to run the python file**

Use WinScp to transfer the python file to the cluster.

**Run the spark submit command:**

**>> spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 datewise_bookings_aggregates_spark.py**



**Command to move the csv file to HDFS**

The data is already loaded in Hdfs by the above pyhon file.

**Screenshot of the file in HDFS**