# Load data from AWS RDS to Hadoop
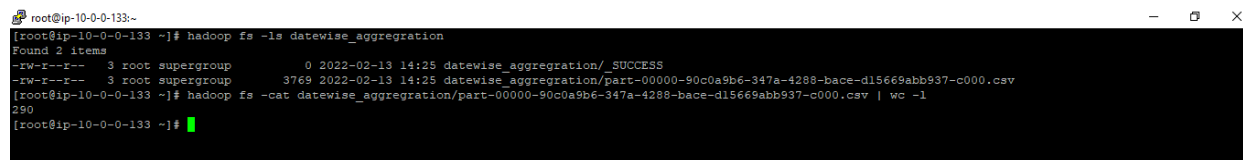
**<Command to run the python file>**

1. Connect to the ec2 instance
2. Switch the user from ec2 to root using. (**sudo –i**)
3. Create a python file with code that can ingest real-time clickstream data from Kafka Server and save it to the local directory. **(vi datewise_bookings_aggregates_spark.py)**
4. Submit the spark job using the python file with spark jar file **(spark2-submit --jars spark-sql-kafka-0-10_2.11-2.3.0.jar datewise_bookings_aggregates_spark.py**)

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

We've specified the path **(/user/root/datewise_aggregration)** for direct saving of the formatted csv and easy retrieval process.

**<Screenshot of the file in HDFS>**