# Load Booking Data Aggregation into Hadoop

1. pyspark file "datewise_bookings_aggregates_spark.py" created to aggregate total number of booking by pickup date

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



2. Verify the aggregated Data

   ```
   hadoop fs -ls /user/root/datewise_bookings_agg/
   ```



   ```
   hadoop fs -cat /user/root/datewise_bookings_agg/part-00000-e993f8e3-a13f-
   4c91-8e55-4f1520f4eaf8-c000.csv | head -n 5
   ```