# Load data from AWS RDS to Hadoop

**Steps to calculate date wise bookings aggregate data**

1. **Establishing spark connection:**

```
spark = SparkSession \
        .builder \
        .appName('aggregateBatchData') \
        .master('yarn') \
        .getOrCreate()
```

2. **To read data from csv file extracted from AWS RDS and stored in HDFS:**

```
df=spark.read.csv('/user/hadoop/bookings-
data/',header=False,inferSchema = True)
```

3. **To add column headers according to given data:**

```
new_columns=["booking_id","customer_id","driver_id","customer_app_version","
customer_phone_os_version","pickup_lat","pickup_lon","drop_lat","drop_lon","
pickup_timestamp","drop_timestamp","trip_fare","tip_amount","currency_code",
"cab_color","cab_registration_no","customer_rating_by_driver","rating_by_cus
tomer","passenger_count"]

new_df = df.toDF(*new_columns)
```

4. **To create a new column with date extracted from pickup_timestamp column:**

```
new_df = new_df.withColumn("date", F.to_date(F.col("pickup_timestamp")))
```

5. **To get the datewise bookings aggregate:**
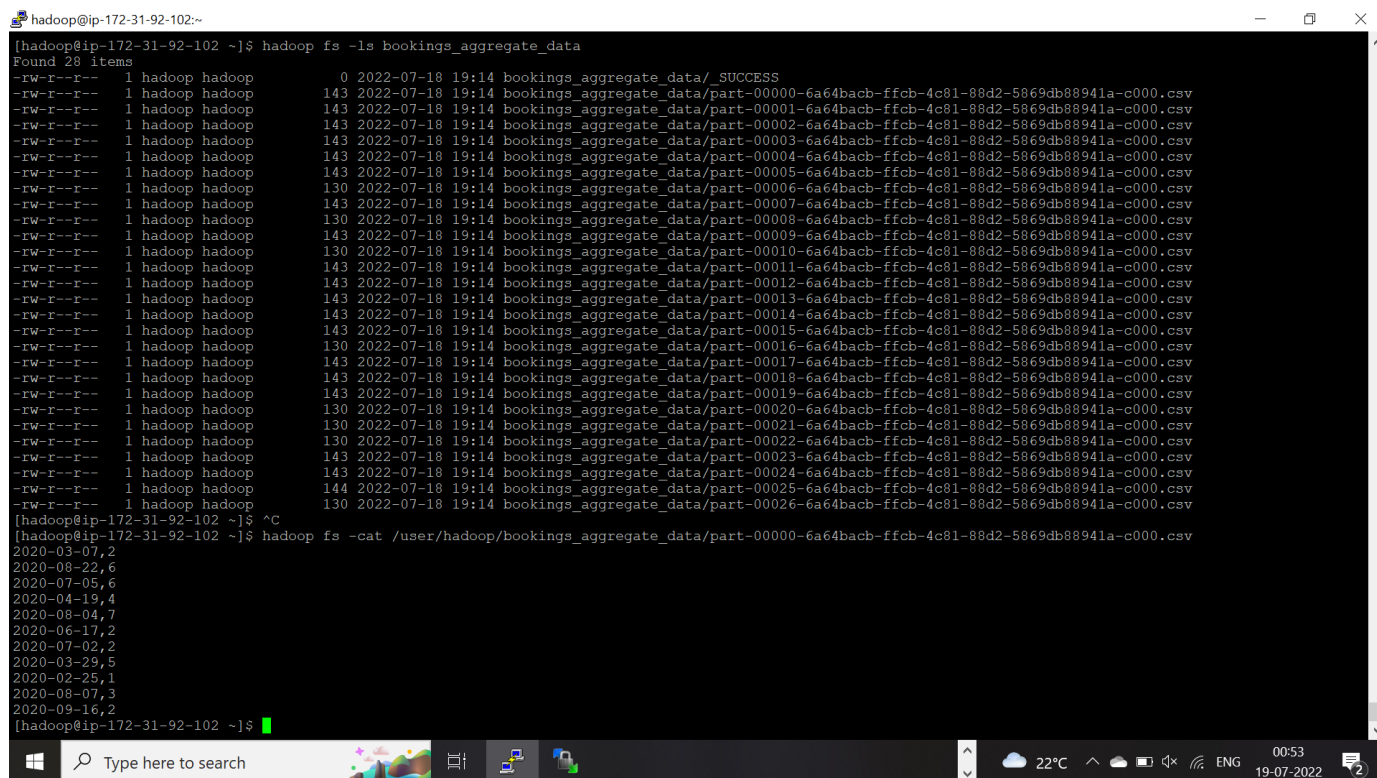
```
aggregate_df = new_df.groupby('date').count()
```

6. **To write the resultant dataframe in csv files in HDFS:**

```
aggregate_df.write.csv('/user/hadoop/bookings_aggregate_data/')
```

7. **Command to run datewise_bookings_aggregates_spark.py file:**

```
spark-submit datewise_bookings_aggregates_spark.py
```

8. **Screenshot of the csv files and data:**