# Aggregate Bookings data and load into Hadoop

**&lt;Command to run the python file&gt;**

*"spark2-submit datewise_bookings_aggregates_spark.py"*

**&lt;Command to move the csv file to HDFS&gt;**

Dataframe is directly stored into the HDFS from the above code. Line 35 does it.
*"aggDF.repartition(1).write.format('com.databricks.spark.csv').save('/user/ec2-user/cab_ride_analysis/aggBookings/results',
header = 'true')"*

**&lt;Screenshot of the file in HDFS&gt;**

https://drive.google.com/file/d/1KtjCnBOtiwFPsJwnTjpte6n0g1BW0bVP/view?usp=sharing

```
[hdfs@ip-10-0-0-52 ~]$ hadoop fs -cat /user/ec2-user/cab_ride_analysis/aggBookings/results/part-00000-51d02429-6d87-4c7f-a98a-5064215a6fa7-c000.csv
Date,Bookings_Count
2020-08-24,1
2020-07-24,5
2020-08-05,3
2020-01-21,1
2020-08-28,5
2020-04-30,6
2020-10-04,5
2020-09-24,6
2020-03-07,2
2020-03-13,2
2020-02-04,4
2020-02-15,3
2020-05-23,4
```