

# Bicycle Sharing Demand

## Bike\_sharing\_project\_file.ipynb

Initial Data Exploration and Transformation is done in data bricks platform and have attached the codebase as Ipython Notebook which explicitly displays each intermediate result and model performances.

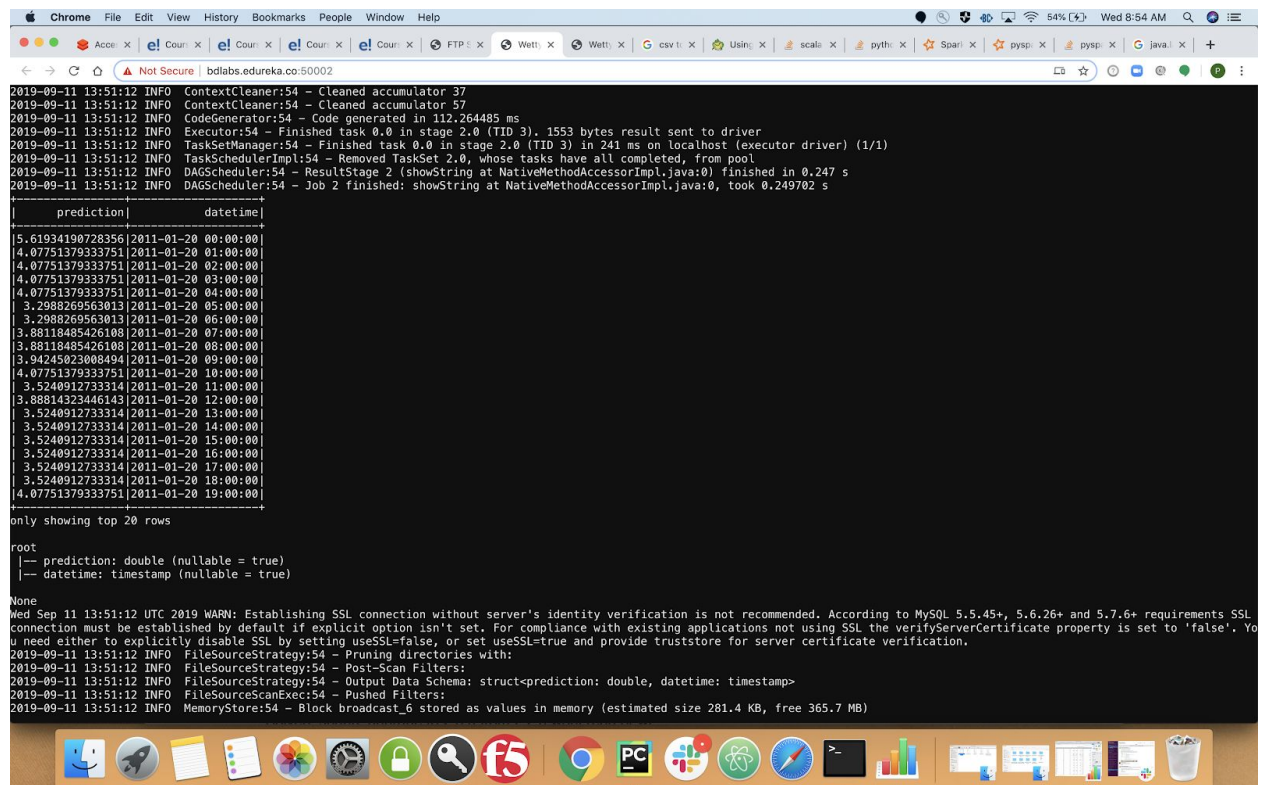
## bike\_sharing\_model\_generation.py

This file is used to give a trained model taking training files as input by cleaning them and using one of the algorithms which gave best results GBRegressor which gave least root mean square error. This pyspark code gives trained model files as **bike\_sharing\_gbt\_file.model**

## bike\_sharing\_prediction.py

This file is uses the model generated out of previous file and predicts the bike sharing demand on the test files given. Then finally outputs the predictions as csv file with name **predictions.csv** and also saves them in RDBMS as shown below.

Predictions are printed below which are in data frame



```
2019-09-11 13:51:12 INFO ContextCleaner:54 - Cleaned accumulator 37
2019-09-11 13:51:12 INFO ContextCleaner:54 - Cleaned accumulator 57
2019-09-11 13:51:12 INFO CodeGenerator:54 - Code generated in 112.264485 ms
2019-09-11 13:51:12 INFO Executor:54 - Finished task 0.0 in stage 2.0 (TID 3). 1553 bytes result sent to driver
2019-09-11 13:51:12 INFO TaskSetManager:54 - Finished task 0.0 in stage 2.0 (TID 3) in 241 ms on localhost (executor driver) (1/1)
2019-09-11 13:51:12 INFO TaskSchedulerImpl:54 - Removed TaskSet 2.0, whose tasks have all completed, from pool
2019-09-11 13:51:12 INFO DAGScheduler:54 - ResultStage 2 (showString at NativeMethodAccessorImpl.java:0) finished in 0.247 s
2019-09-11 13:51:12 INFO DAGScheduler:54 - Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 0.249702 s

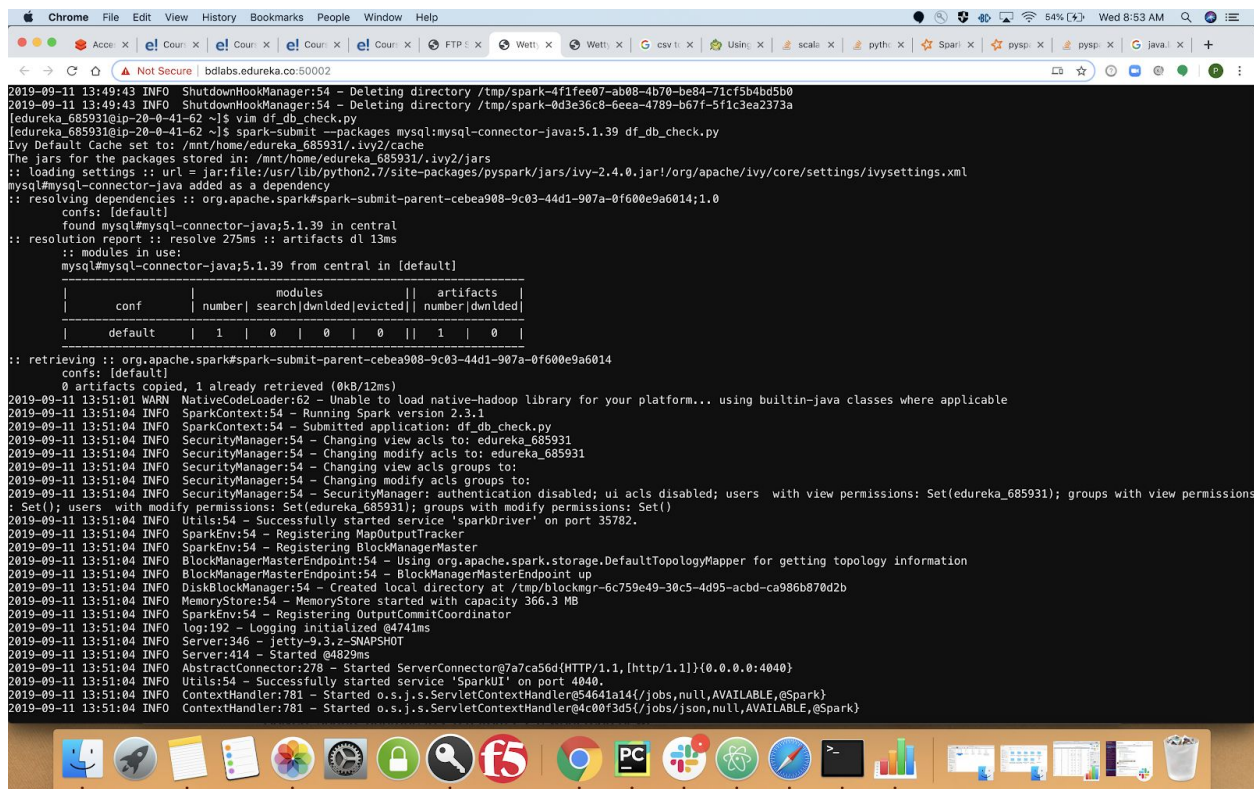
+-----+-----+
| prediction | datetime |
+-----+-----+
| 5.61934190728356 | 2011-01-20 00:00:00 |
| 4.07751379333751 | 2011-01-20 01:00:00 |
| 4.07751379333751 | 2011-01-20 02:00:00 |
| 4.07751379333751 | 2011-01-20 03:00:00 |
| 4.07751379333751 | 2011-01-20 04:00:00 |
| 3.2988269563013 | 2011-01-20 05:00:00 |
| 3.2988269563013 | 2011-01-20 06:00:00 |
| 3.88118485426108 | 2011-01-20 07:00:00 |
| 3.88118485426108 | 2011-01-20 08:00:00 |
| 3.94245823008494 | 2011-01-20 09:00:00 |
| 4.07751379333751 | 2011-01-20 10:00:00 |
| 3.5240912733314 | 2011-01-20 11:00:00 |
| 3.88814323446143 | 2011-01-20 12:00:00 |
| 3.5240912733314 | 2011-01-20 13:00:00 |
| 3.5240912733314 | 2011-01-20 14:00:00 |
| 3.5240912733314 | 2011-01-20 15:00:00 |
| 3.5240912733314 | 2011-01-20 16:00:00 |
| 3.5240912733314 | 2011-01-20 17:00:00 |
| 3.5240912733314 | 2011-01-20 18:00:00 |
| 4.07751379333751 | 2011-01-20 19:00:00 |
+-----+-----+

only showing top 20 rows

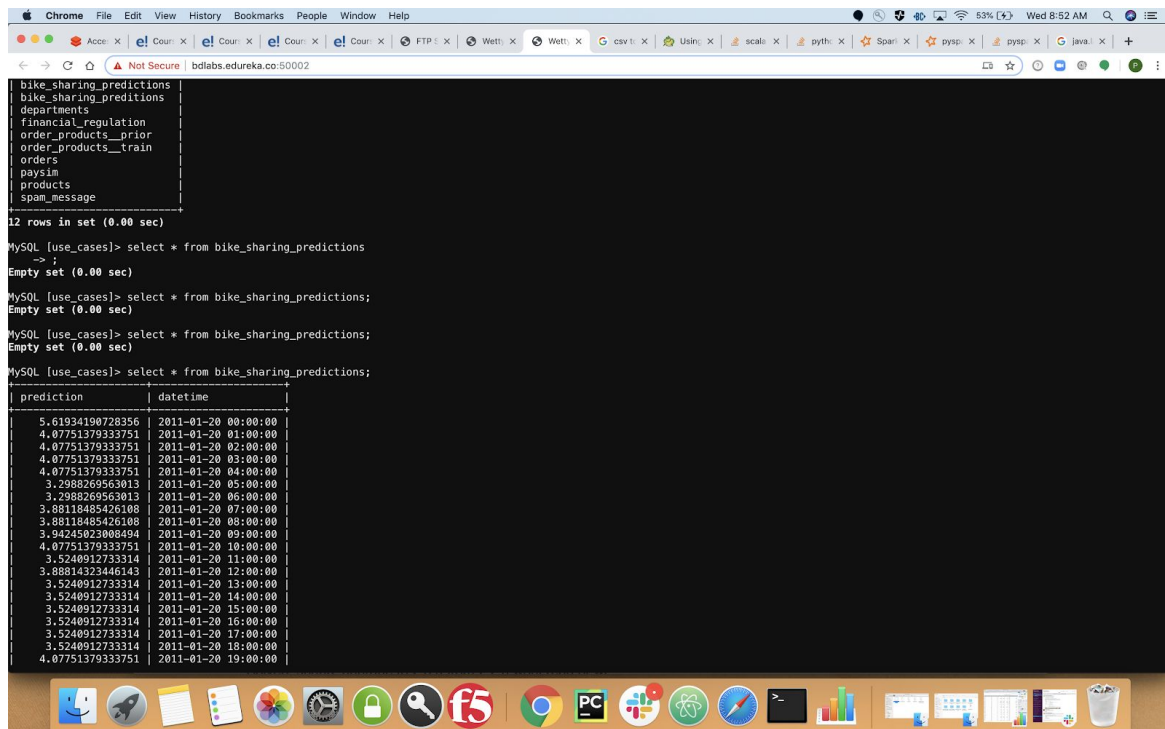
root
|-- prediction: double (nullable = true)
|-- datetime: timestamp (nullable = true)

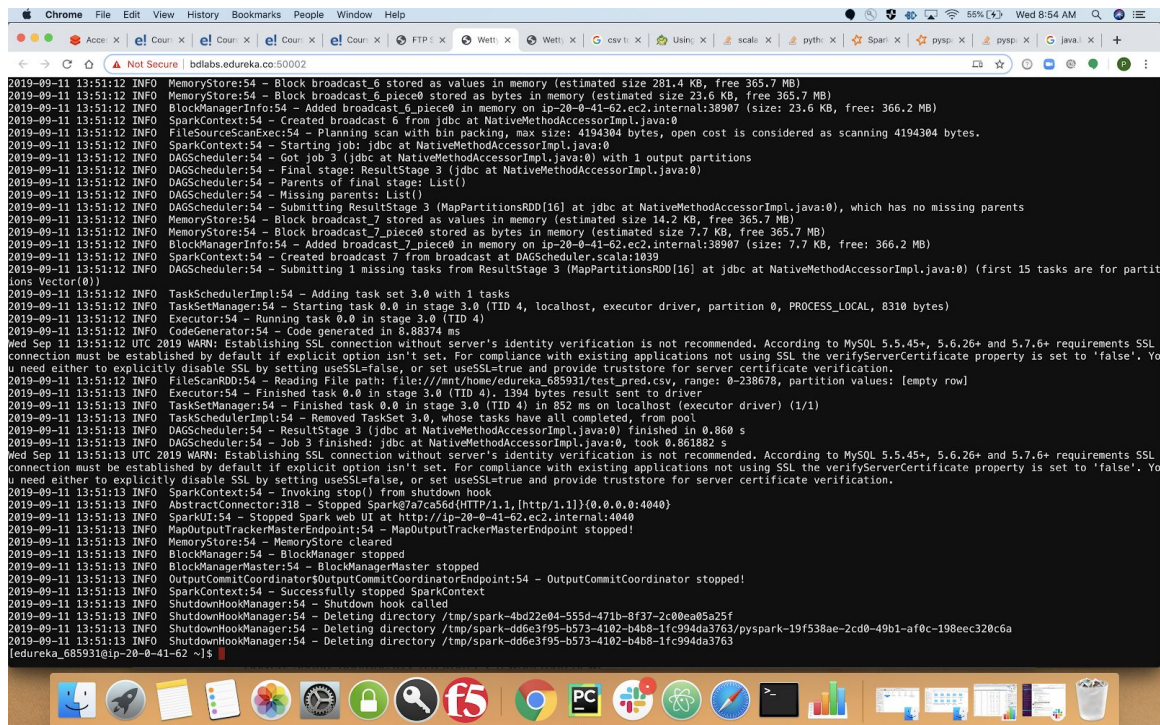
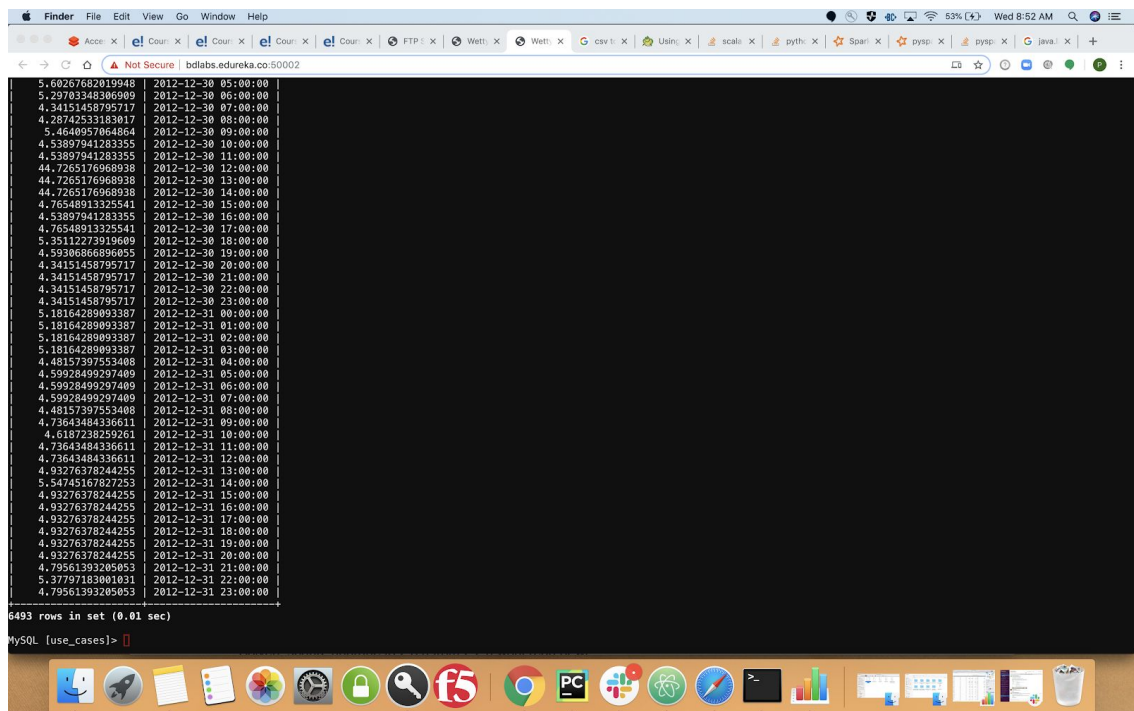
None
Wed Sep 11 13:51:12 UTC 2019 WARN: Establishing SSL connection without server's identity verification is not recommended. According to MySQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection must be established by default if explicit option isn't set. For compliance with existing applications not using SSL the verifyServerCertificate property is set to 'false'. You need either to explicitly disable SSL by setting useSSL=false, or set useSSL=true and provide truststore for server certificate verification.
2019-09-11 13:51:12 INFO FileSourceStrategy:54 - Pruning directories with:
2019-09-11 13:51:12 INFO FileSourceStrategy:54 - Post-Scan Filters:
2019-09-11 13:51:12 INFO FileSourceStrategy:54 - Output Data Schema: struct<prediction: double, datetime: timestamp>
2019-09-11 13:51:12 INFO FileSourceScanExec:54 - Pushed Filters:
2019-09-11 13:51:12 INFO MemoryStore:54 - Block broadcast_6 stored as values in memory (estimated size 281.4 KB, free 365.7 MB)
```

Predictions are saved in RDBMs tables to run this pyspark code required driver packages are also passed as parameters.



The predictions along with the datetime are pushed to RDBMS from the dataframe





## Streaming\_prediction.scala

The streaming prediction is done in scala it uses existing model generated and kafka to stream the features and predict the demand then finally push them to RDBMS connected.