-Sricharan Pamuru (20)

# 2. Spark ML Lib Application

The below screenshot depicts the training data provided in terms of three constraints, gas, bars and restaurants. Prediction is found in the console and it can be visualized as follows:

Prediction output is displayed and the bars are obtained as output based on data collected. Comparison with training data brings a combined result.