# ASSIGNMENT -5 (Kafka + Spark)

Sruthy Suji
23AD140
AI-DS C

## 1) Create a new topic – greentaxi-topic



## 2) Stream Data @ 0.1 Sec Interval using Kafka_Taxi_Producer.py

3) Show Top 5 Trending Pickup Local – "PULocationID" [5] and Drop Location – "DOLocationID" [6]