

## DA 331 - Big Data Analytics : Tools & Techniques

### Lab 6

**Instructor:** Dr. Chiranjib Sur

Dataset link: <https://www.kaggle.com/datasets/elemento/nyc-yellow-taxi-trip-data>

Given the dataset, find the following:

1. You want the biggest (average) trips, which location will you choose?
2. You want the biggest (average) trips, which location will you choose? (in pandas) Find the difference in running time.
3. Find the locations where the maximum number of passengers arrives.
4. Find the locations where the maximum number of passengers starts.
5. Find the locations where the maximum number of passengers starts in a day.

You may need this to process the dates.

```
from pyspark.sql.functions import year, month, day  
  
extracted df = df.select(year("date").alias("year"), /  
month("date").alias("month"), day("date").alias("day"))  
  
extracted df.show()
```