

# AWS Big Data Project

## Problem statement:

Imagine that you are working as an analyst in a famous taxi app company. Your organization provides hassle-free travel to people all around the world. You are provided with an AWS account to perform certain analytical tasks using AWS cloud services.

Link for the dataset and resources:

<https://intellipaate-course-attachments.s3.ap-south-1.amazonaws.com/Hadoop/Hadoop+Datasets-20200609T120700Z-001.zip>

## Dataset description:

In here, you have a predefined dataset (**yellow.csv**), having more than 15 columns.

The dataset has different attributes as follows:

vendor\_id string,  
pickup\_datetime string,  
dropoff\_datetime string,  
passenger\_count int,  
trip\_distance DECIMAL(9,6),  
pickup\_longitude DECIMAL(9,6),  
pickup\_latitude DECIMAL(9,6),  
rate\_code int,  
store\_and\_fwd\_flag string,  
dropoff\_longitude DECIMAL(9,6),  
dropoff\_latitude DECIMAL(9,6),  
payment\_type string,

fare\_amount DECIMAL(9,6),  
extra DECIMAL(9,6),  
mta\_tax DECIMAL(9,6),  
tip\_amount DECIMAL(9,6),  
tolls\_amount DECIMAL(9,6),  
total\_amount DECIMAL(9,6),  
trip\_time\_in\_secs int

## Tasks:

- Stream the taxi data to a Kinesis Stream
- Link the Kinesis Stream to a Kinesis Firehose delivery stream
- Deliver the streamed data to Amazon S3
- Visualize the relationship between the total fare and the distance covered in trips
- Visualize the percentage of payment types using a pie chart
- Visualize the relationship between the passenger count and the tip amount

## Use EMR to find out the following details:

1. What is the total number of trips?
2. What is the total revenue generated by all the trips? (The fare is stored in the column total\_amount)
3. What fraction of the total is paid for tolls? (Toll fares are stored in tolls\_amount)
4. What fraction of the total is paid as driver tips? (Tips are stored in tip\_amount)
5. What is the average trip amount?
6. What is the average distance of the trips? (Distances are stored in the column trip\_distance)
7. How many different payment types are used?
8. For each payment type, display the following details:
  - Average fare generated
  - Average tip
  - Average tax (Tax is stored in column mta\_tax)
9. On average, which hour of the day generates the highest revenue?