# DATA CURATION PROJECT: PRELIMINARY OBSERVATIONS OF THE YELP DATASET

## DATASET NAME

Yelp Dataset (from Kaggle) - https://www.kaggle.com/datasets/yelp-dataset/yelp-dataset/data

## DATA UNDERSTANDING AND QUALITY ISSUES:

o PROVIDE DESCRIPTIVE STATISTICS FOR NUMERIC ATTRIBUTES.
o IDENTIFY MISSING VALUES.
o CHECK FOR DUPLICATES.
o IDENTIFY ANY ANOMALIES IN THE DATA (BUSINESSES WITH AN EXTREMELY HIGH OR LOW NUMBER OF REVIEWS).

## DATA CLEANING:

o HANDLE MISSING VALUES (IMPUTE OR REMOVE).
o HANDLE DUPLICATES (REMOVE).
o NORMALIZE TEXT ATTRIBUTES (CONVERT ALL TEXT TO LOWERCASE).

## DATA INTEGRATION:

o INTEGRATE TWO DATASETS (YELP REVIEWS) BASED ON COMMON ATTRIBUTES (**BUSINESS_ID**).
o PERFORM SCHEMA MATCHING AND MAPPING IF REQUIRED.

## DATA TRANSFORMATION (ETL PROCESSES):

o EXTRACT RELEVANT ATTRIBUTES FOR ANALYSIS.
o TRANSFORM ATTRIBUTES (EXTRACT THE PRIMARY CATEGORY FROM THE **CATEGORIES** COLUMN).
o LOAD INTO A STRUCTURED FORM SUITABLE FOR QUERYING OR ANALYTICS.

## DATA ANALYSIS:

o IDENTIFY THE TOP-RATED BUSINESSES.
o FIND THE AVERAGE RATING PER CITY OR STATE.
o ANALYZE THE DISTRIBUTION OF BUSINESSES ACROSS DIFFERENT CATEGORIES.
o EXPLORE THE RELATIONSHIP BETWEEN THE NUMBER OF REVIEWS AND RATINGS.

## DATA PROVENANCE:

o DOCUMENT THE SOURCE OF THE DATA.
o TRACK ANY CHANGES OR TRANSFORMATIONS MADE TO THE DATA.
o STORE METADATA OR INFORMATION ABOUT DATA PROCESSING STEPS.

## VISUALIZATION:

o PLOT THE DISTRIBUTION OF RATINGS.
o VISUALIZE THE NUMBER OF BUSINESSES IN DIFFERENT CATEGORIES.
o MAP THE BUSINESSES BASED ON THEIR LATITUDE AND LONGITUDE TO VISUALIZE THEIR DISTRIBUTION.

Other/Alternative dataset - https://github.com/bookingcom/ml-dataset-mdt [Booking.com]