Lab 1

Datasets

1 Electric Vehicle population Dataset

problem statement: The dataset contains information about Battery Electric Vehicles and plug-in Hybrid Electric Vehicles registered in Washington state. it includes details like vehicle type, make, model, Electric range and location. The dataset has 223,995 rows and 17 columns with missing values. possible problem statements include predicting EV adoption based on vehicle features and geographic location, analyzing regional factors influencing EV popularity

Number of rows = 223995
Number of columns = 17

Columns

① VIN
② county
③ City
④ State
⑤ postal Code
⑥ Model Year
⑦ Make
⑧ Model
⑨ Electric Vehicle Type
⑩ Clean Alternative fuel Vehicle Eligible
⑪ Electric Range
⑫ Base MSRP
⑬ Legislative District
⑭ DOL Vehicle ID
⑮ vehicle location

⑯ Electric Utility
⑰ 2020 census Tract

## 2. Laptop price Dataset

Problem statement: This dataset contains 4768 records of laptops with various hardware specifications and prices it is designed for predictive modeling, price estimation and exploratory data analysis. this dataset includes real world values to reflect the factors that influence laptop prices. This dataset can be used for pr. Building machine learning models to predict laptop price based on specification, Market Analysis & Feature importance study.

number of rows = 4769
number of columns = 11

### Coloums

1. Brand
2. processor
3. RAM
4. Storage
5. GPU
6. Screen Size
7. Resolution
8. Battery life
9. weight
10. operating System
11. price.

3 Indian Bike Sales Dataset

Problem statement : This dataset contains records of motorcycle sales across various indian state, covering top brands like Honda, Royal Enfield, TVS, Yamaha, Hero, Bajaj, KTM and Kawasaki. The dataset includes key attributes such as average daily distance traveled, engine capacity, fuel type etc. It provides insights into bike sales trends, market demand, and resale value across different city tiers

number of rows = 10001
number of columns = 15

Colowns

① State
② Avg Daily Distance
③ Brand
④ Model
⑤ Price
⑥ Year of Manufacture
⑦ Engine Capacity
⑧ Fuel Type
⑨ Mileage
⑩ owner Typ
⑪ Registration
⑫ insurance
⑬ seller Type
⑭ Resale price
⑮ city Tier

(4) Road Accident Survival Dataset

problem Statement:
This dataset contains detailed records of simulated road accident data, focusing on factors influencing survival outcomes. The dataset includes demographic, behavioural and situational attributes, providing valuable insight into how various factor impact the survival probability during road accidents.

number of rows = 201
number of colourmms = 6

colourmms
① Age
② Gender
③ Speed-of-impact
④ Helmet-Used
⑤ Seatbelt-Used
⑥ Survived

(5) phone usage in India

problem Statement:
The dataset pre represents simulated phone usage data for Indian users. It contains 17,886 rows and 13 columns, each row representing an individual's phone usage details. The data reflects various aspects of mobile phone behaviour including demographics, phone brand and app usage patterns

number of rows = 14687
number of columns = 16

① User ID
② Age
③ Gender
④ Location
⑤ Phone Brand
⑥ OS
⑦ Screen Time
⑧ Data Usage
⑨ Calls Duration
⑩ Number of Apps install
⑪ Social Media Time
⑫ E-commerce Spend
⑬ Streaming Time
⑭ Gaming Time
⑮ Monthly Recharge
⑯ primary useage

3/3/2025