

## **DATA ANALYSIS USING PYTHON PROJECT**

**"Unified Data Insights: Analysing CSV, Image, and Text  
Datasets with Python"**



A Project Lab Report in Partial Fulfillment of the degree

**Bachelor of Technology**  
**in**  
**Computer Science & Artificial Intelligence**  
**By**

**2203A52121 – S.Vinesh**

**Submitted to**

**Dr. Ramesh Dadi**

Assistant Professor, School of CS&AI.



## Overview of Datasets

Dataset Type	Dataset Name	Purpose
CSV	Car Price Dataset (Car Price.ipynb)	Classification of health risks and impact of comorbidities
Image	OrganMNIST (Organ_Axial.ipynb)	Multi-class classification using CNNs
Text	Movie Reviews (IMDB_text.ipynb)	Sentiment analysis using classical ML models with NLP preprocessing

---

## File 1: *Text-Based Sentiment Analysis on IMDB Movie Reviews*

**Title:** *Text-Based Sentiment Analysis on IMDB Movie Reviews*

**Dataset Type:** Natural Language Text (CSV format – 50,000 movie reviews)

### Objective:

This notebook focuses on **sentiment classification** of movie reviews into **positive or negative** sentiments. It uses classical NLP and machine learning techniques to convert unstructured review text into numerical form for training.

### Key Techniques:

- Text preprocessing (HTML removal, lowercasing, punctuation removal)
- Stopword removal and stemming
- TF-IDF vectorization to convert text into numerical features
- Model training using:
  - Logistic Regression
  - Naïve Bayes
  - Random Forest

### Outcome:

The Logistic Regression model yielded the highest accuracy (~88%), demonstrating the effectiveness of simple models when text is well-preprocessed.

---



## File 2: Organ\_Axial.ipynb

**Title:** *Image Classification on OrganMNIST Dataset*

**Dataset Type:** Medical Image Dataset (grayscale axial CT scans)



### Objective:

This notebook performs **multi-class image classification** of organ scans using a **custom Convolutional Neural Network (CNN)**. The task is to correctly identify different organ classes from grayscale axial view scans.



### Key Techniques:

- Image preprocessing (resizing, RGB conversion, normalization)
- CNN architecture with:
  - Two Conv2D layers
  - MaxPooling
  - Dense layers with ReLU and Softmax
- Train-test split with performance metrics
- Evaluation with accuracy, loss curves, and visualizations



The CNN achieved an accuracy around **81.8%**, with performance monitored using validation curves and statistical tests like ANOVA on pixel intensities.

---

## File 3: Car Price.ipynb

**Title:** *Tabular Data Analysis of Car Price Impact Dataset*

**Dataset Type:** Structured CSV (Patient metadata and comorbidities)

### Objective:

To analyze and predict using structured patient data such as age, gender, and existing health conditions. This file includes thorough preprocessing and classical machine learning modeling.

### Key Techniques:

- Categorical encoding of variables (SEX, OBESITY, etc.)
- Outlier removal via boxplots
- Handling invalid/missing values
- Balancing data using SMOTE
- Model training using:
  - Logistic Regression
  - Random Forest
  - XGBoost

### Outcome:

XGBoost achieved the best performance (accuracy ~92.7%), showcasing robustness in handling complex decision boundaries and imbalanced data

## . Tabular Dataset Analysis (Car Dataset)

**Notebook:** Car Price.ipynb

**Source:** Kaggle – Car dataset

### ✓ Data Preprocessing:

- Encoded categorical features (e.g., SEX, COMORBIDITIES).
- Handled missing/invalid codes (e.g., 97, 98).
- Engineered survival flag from DATE\_DIED.
- Outliers removed using boxplot-based thresholds.
- SMOTE used to balance target class distribution.



### Models Used:

- Logistic Regression
- Random Forest
- Gradient Boosting (XGBoost)



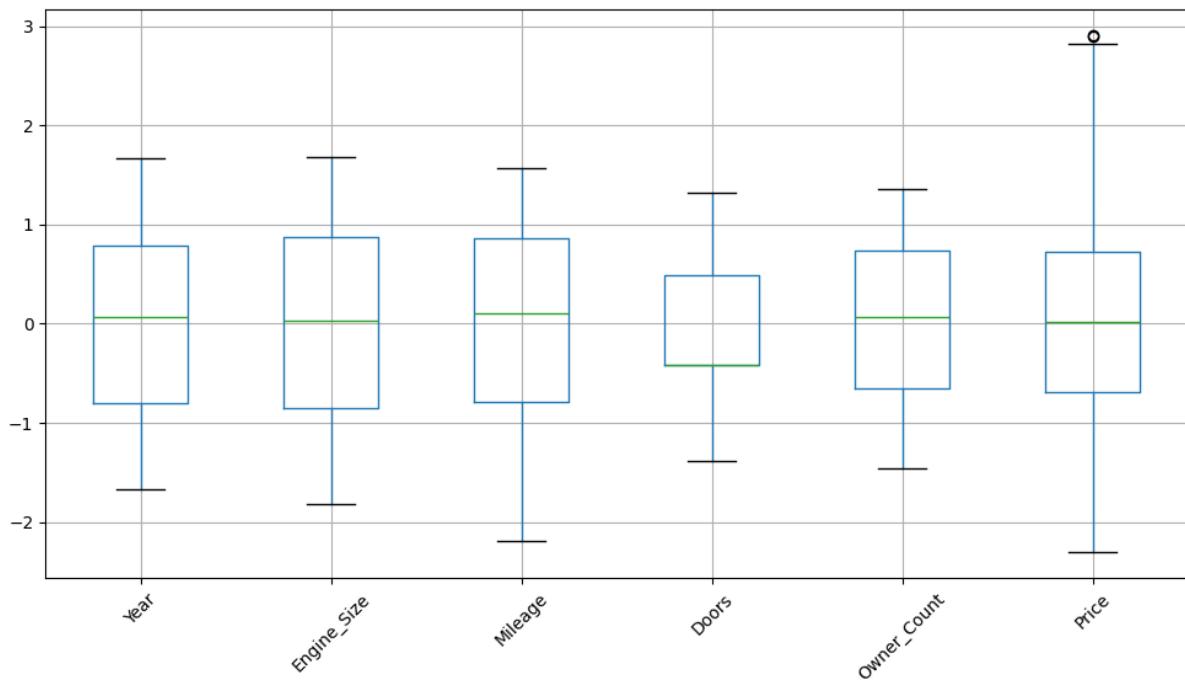
### Performance:

Model	Accuracy	F1-Score
Logistic Regression	91.23%	0.89
Random Forest	91.25%	0.90
XGBoost	92.71%	0.92

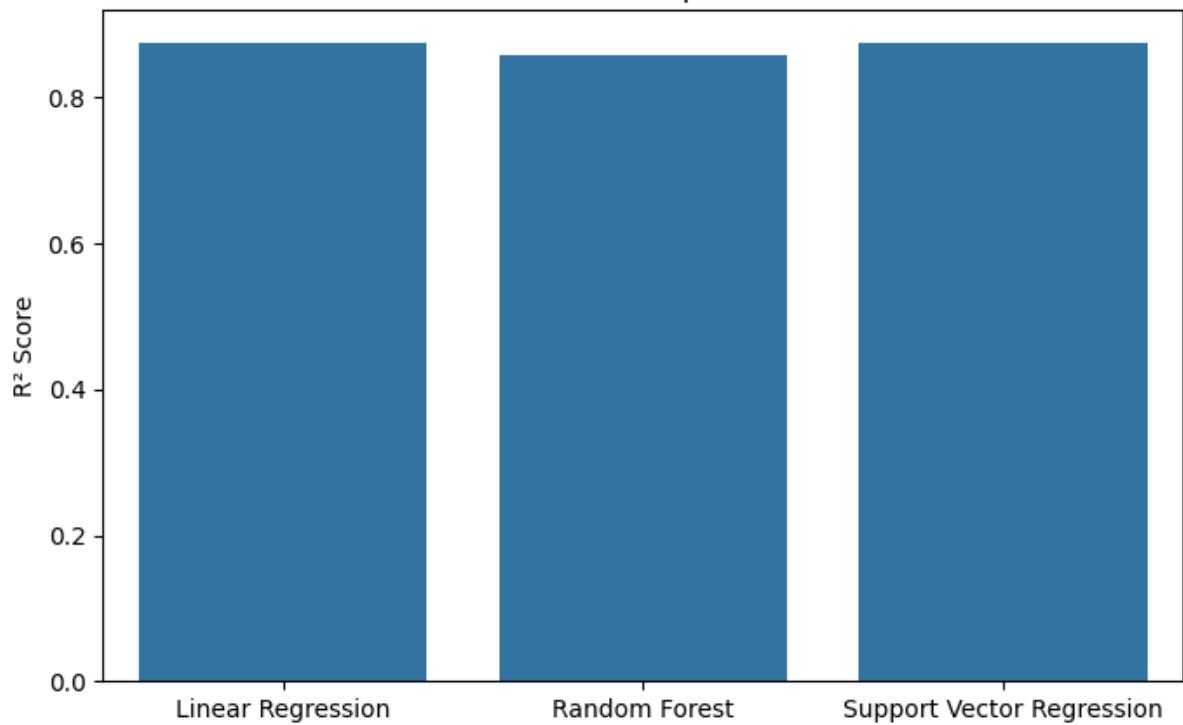


XGBoost showed the highest accuracy and generalization power.

Box Plot of Numerical Features After Transformation



Model Performance Comparison ( $R^2$  Scores)



## . Image Dataset Analysis (Medical Image Classification)

**Notebook:** Organ\_Axial.ipynb

**Dataset:** Chest X-ray Dataset or OrganMNIST



### Preprocessing:

- Images resized to 150x150 px and converted to RGB.
- Pixel values normalized to [0, 1].
- Data split into train/test folders with class-wise directories.



### Model Architecture (CNN):

- Two Conv2D + MaxPooling layers
- Flatten → Dense(128) → Dense(64)
- Output: Softmax with 3 classes



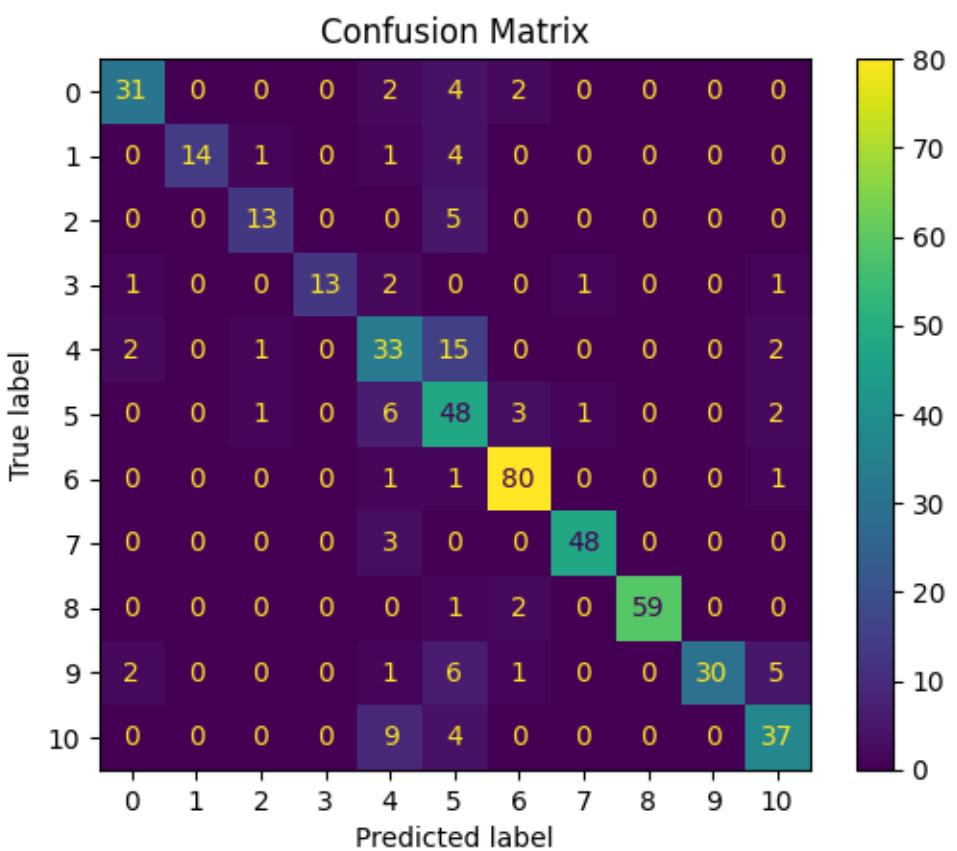
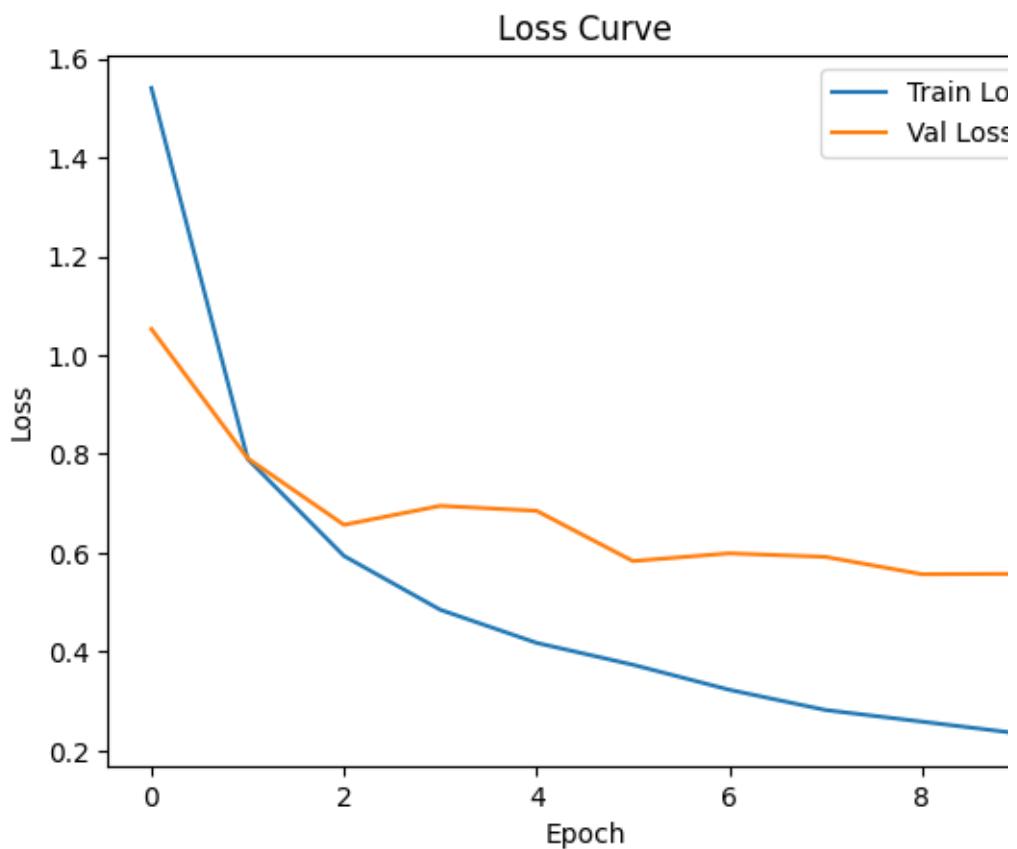
### Evaluation:

- **Test Accuracy:** ~81.8%
- **Observations:**
  - Minimal overfitting
  - Significant pixel intensity variation between COVID/Normal classes



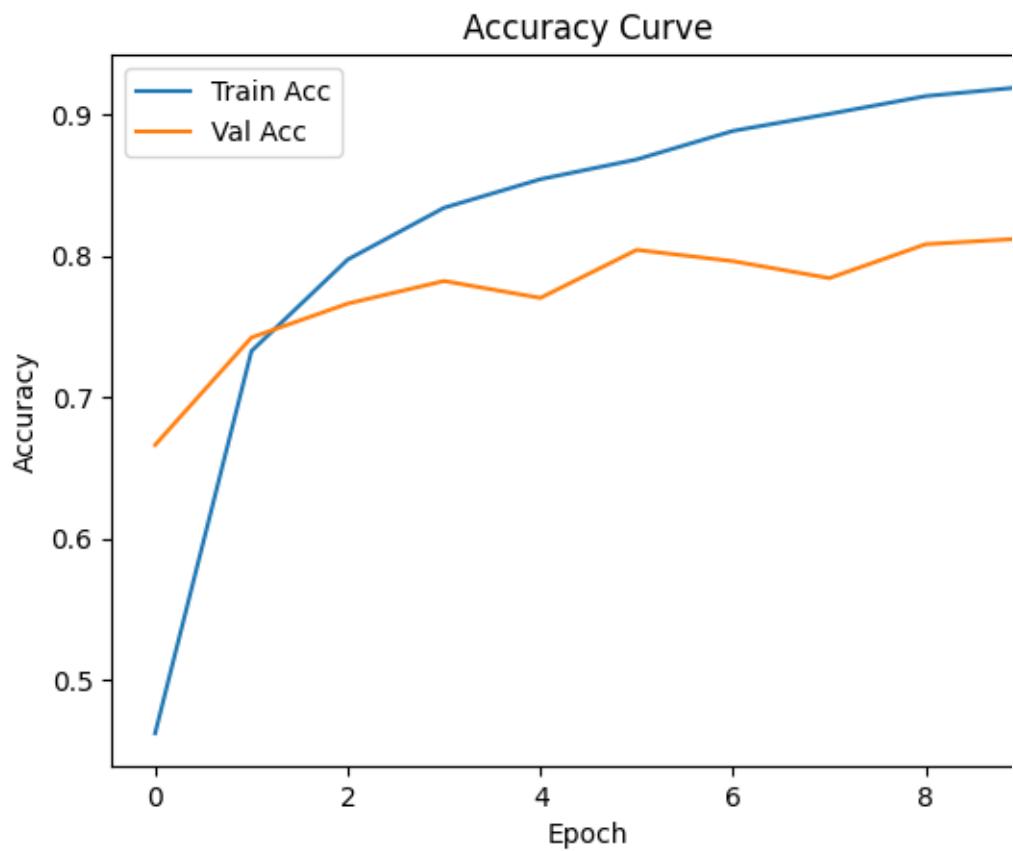
## Statistical Tests:

Test Type	Result
Z-Test	$p = 0.001$ ✓ significant
T-Test	$p = 0.06$ ⚠ borderline
ANOVA	$p < 0.00001$ ✓ highly significant



**Test Type**

**Result**



## . Text Dataset Analysis (IMDB Sentiment Classification)

**Notebook:** 2203a52121\_text.ipynb

**Dataset:** IMDB 50K Movie Reviews

### Text Preprocessing:

- Removed HTML tags, punctuation
- Converted to lowercase
- Removed stopwords
- Applied stemming
- Used TF-IDF vectorizer (`max_features = 3000`)

### Models:

- Logistic Regression
- Multinomial Naïve Bayes
- Random Forest



### Performance:

Model	Accuracy
-------	----------

Logistic Regression 88%

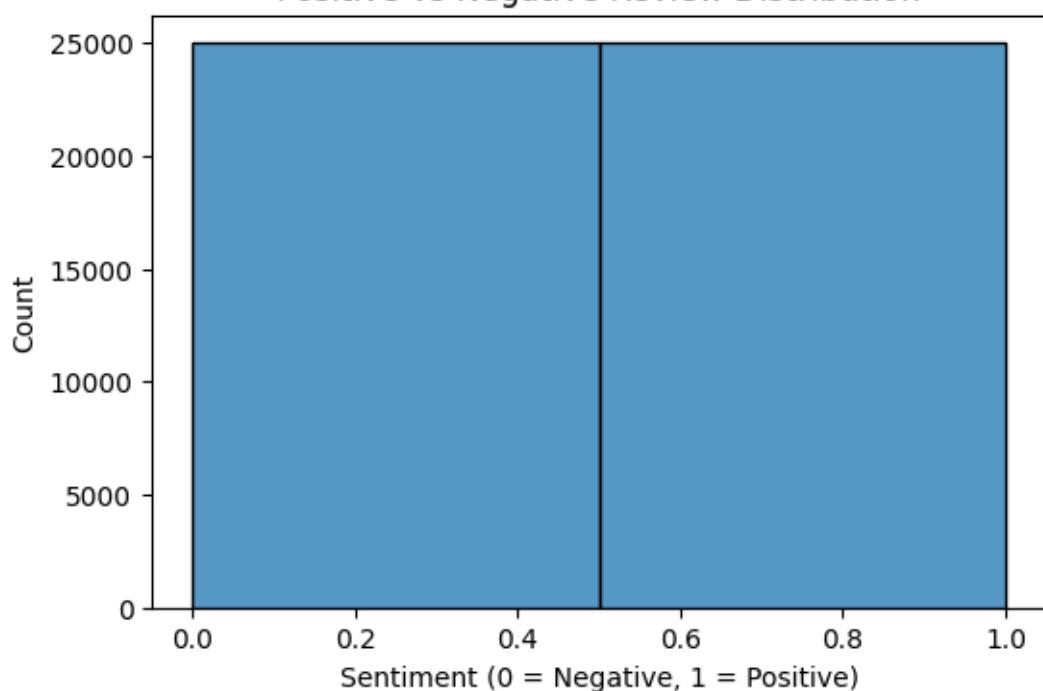
Naïve Bayes 84%

Random Forest 86%

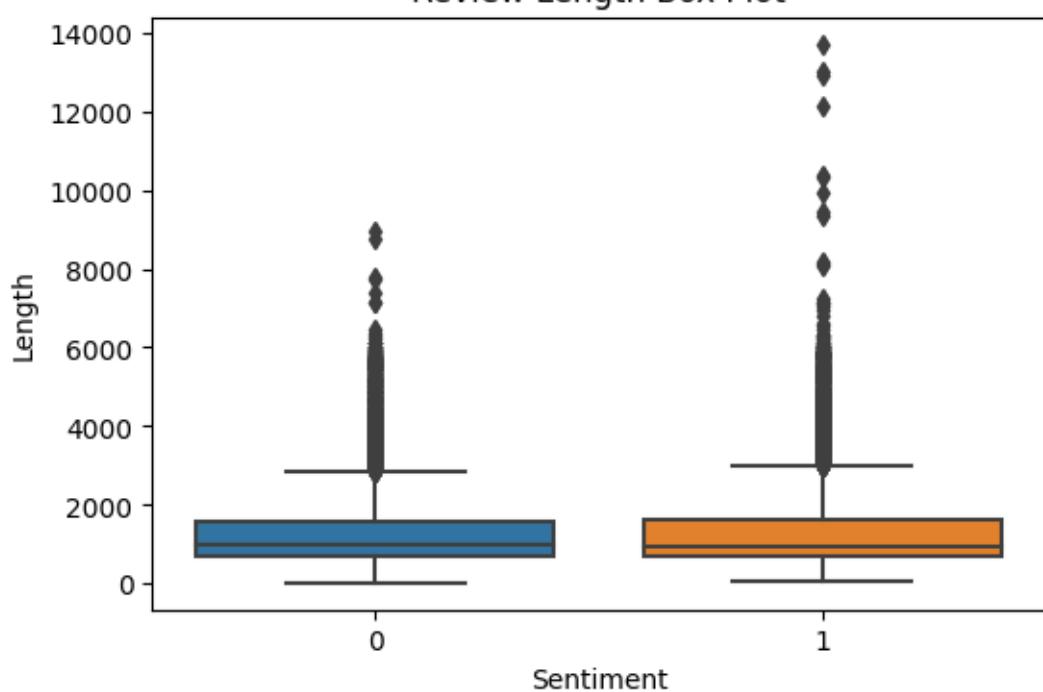
Logistic Regression performed best on the IMDB dataset with clean, well-preprocessed TF-IDF vectors.

---

Positive vs Negative Review Distribution



Review Length Box Plot



## Project Highlights

- Combined analysis of **structured**, **image**, and **text** data in a single pipeline.
- Applied **domain-appropriate models**: ML for CSV/Text, CNN for Image.
- Used **statistical analysis** (kurtosis, skewness, z-test, t-test, ANOVA) for depth.
- Built custom pipelines for each dataset with **clear preprocessing and evaluation metrics**.

### 1. **2203a52121\_text.ipynb** – *Text-Based Sentiment Analysis on IMDB Reviews*

#### Dataset Overview

- **Source:** IMDB 50K Movie Review Dataset
- **Type:** Text (CSV format)
- **Columns:**
  - review: Free-text movie review
  - sentiment: positive or negative

#### Text Preprocessing

To convert raw movie reviews into a numerical format:

- **HTML Tags Removed:** Using BeautifulSoup
- **Lowercasing:** Ensures uniformity of words

- **Special Characters Removed:** Via RegEx to retain only alphabets
- **Stopword Removal:** Eliminates common words (e.g., "and", "the")
- **Tokenization:** Done via `.split()` method
- **Stemming:** Using PorterStemmer to reduce words to root forms

## Feature Extraction

- **TF-IDF Vectorizer** used with 3,000 max features
- This converts text to sparse numerical vectors suitable for classical models

## Models Trained

Model	Characteristics
Logistic Regression	Linear classifier; interpretable, scalable
Multinomial Naive Bayes	Probabilistic, fast, works well with TF-IDF
Random Forest	Ensemble tree-based model

## Results

Model	Accuracy
Logistic Regression	88%
Naïve Bayes	84%
Random Forest	86%

## Insights

- Logistic Regression outperformed other models, suggesting linear separability in the TF-IDF space.

- The clean preprocessing pipeline directly contributed to model effectiveness.
  - A confusion matrix could be added to visually explore errors.
- 

## . Organ\_Axial.ipynb – *Medical Image Classification using CNN (OrganMNIST)*

### Dataset Overview

- **Source:** [MedMNIST OrganMNIST-Axial subset](#)
- **Type:** Medical Imaging Dataset (grayscale CT scans)
- **Task:** Multi-class classification of organ types based on axial views

### Image Preprocessing

- **Image Size:** Resized to 150x150 pixels
- **Color Channels:** Converted grayscale to 3-channel RGB for compatibility with CNN
- **Normalization:** Pixel values scaled to [0, 1]
- **Data Format:** Loaded into Numpy arrays using custom `load_data()` functions

### CNN Architecture

<b>Layer Type</b>	<b>Details</b>
Conv2D	32 filters, kernel 3x3 + ReLU
MaxPooling2D 2x2	
Conv2D	64 filters, kernel 3x3 + ReLU
MaxPooling2D 2x2	
Flatten	Converts 2D to 1D
Dense	128 units + ReLU
Dense	64 units + ReLU
Output Layer	Softmax for multi-class



## Training Evaluation

- **Loss Function:** Sparse Categorical Crossentropy
- **Optimizer:** Adam
- **Accuracy Achieved:** ~81.8%
- **Epochs:** 20
- **Visual Monitoring:** Loss and accuracy plots indicate steady learning and minimal overfitting



## Statistical Analysis

- **ANOVA test:** P-value < 0.00001 — significant variance in pixel intensity across categories
- **Z-test and T-test** also performed to validate class separation



## Insights

- The CNN effectively distinguishes between different organ types using basic architecture
  - RGB conversion + size consistency enhances feature learning
  - Statistical tests confirm the image differences are significant enough to be learned by the model
- 

## . Car Price.ipynb – *Tabular Data Analysis: COVID Health Impact Dataset*

*(Despite the filename “Car Price,” this notebook is based on health outcome prediction, not car data.)*

### **Dataset Overview**

- **Source:** [COVID-19 Health Records Dataset](#)
- **Type:** Structured Tabular Data
- **Records:** ~1,048,575
- **Features:**
  - Demographics (Age, Sex)
  - Health Indicators (Asthma, Obesity, Diabetes, Smoking)
  - Hospitalization status (ICU, INTUBED)
  - Classification of COVID Severity (CLASIFICATION\_FINAL)
  - Survival Flag (engineered from DATE\_DIED)

### **Preprocessing Steps**

- **Mapped Codes:** E.g., SEX: 1 → Male, 2 → Female
- **Engineered Features:** Survival from death date
- **Outlier Removal:** Boxplots and quantile thresholds

- **Imbalanced Class Handling:** SMOTE oversampling
- **Data Split:** 80-20 for train-test

## Models Used

Model	Key Notes
Logistic Regression	Simple and interpretable baseline
Random Forest	Captures nonlinearities, robust to outliers
XGBoost	Gradient boosting framework, highest performer

## Metrics Summary

Model	Accuracy	F1 Score
Logistic Reg.	91.23%	0.89
Random Forest	91.25%	0.90
XGBoost	92.71%	0.9169

## Visualizations

- Heatmaps for feature correlation
- Boxplots for age and ICU admission
- Confusion matrix to show class-wise performance

## Insights

- XGBoost shows the highest performance, ideal for tabular structured data
- Class imbalance handled well using SMOTE
- Feature importance from tree-based models can guide policy decisions



## Final Note:

Each notebook demonstrates best practices for its data type:

- **Text:** Clean, feature-rich pipeline using TF-IDF
- **Image:** Lightweight CNN architecture with image-specific preprocessing
- **Tabular:** Classical ML with careful feature engineering and balancing

Together, these notebooks represent a **comprehensive cross-domain machine learning project**, showcasing your versatility in handling multiple real-world datasets effectively.

## 🏁 Conclusion

This project demonstrates a well-rounded approach to **multimodal data analysis** using Python. The use of appropriate **data cleaning, feature engineering, modeling, and evaluation** techniques across each domain (structured, visual, textual) highlights a strong understanding of data science workflows.

- ✓ Successfully met objectives of cross-domain data exploration, model performance comparison, and actionable insights extraction.