

# MACHINE LEARNING

Group Project



# CONTENT

1  
2

3

GITHUB LINK

NUMERICAL DATASET

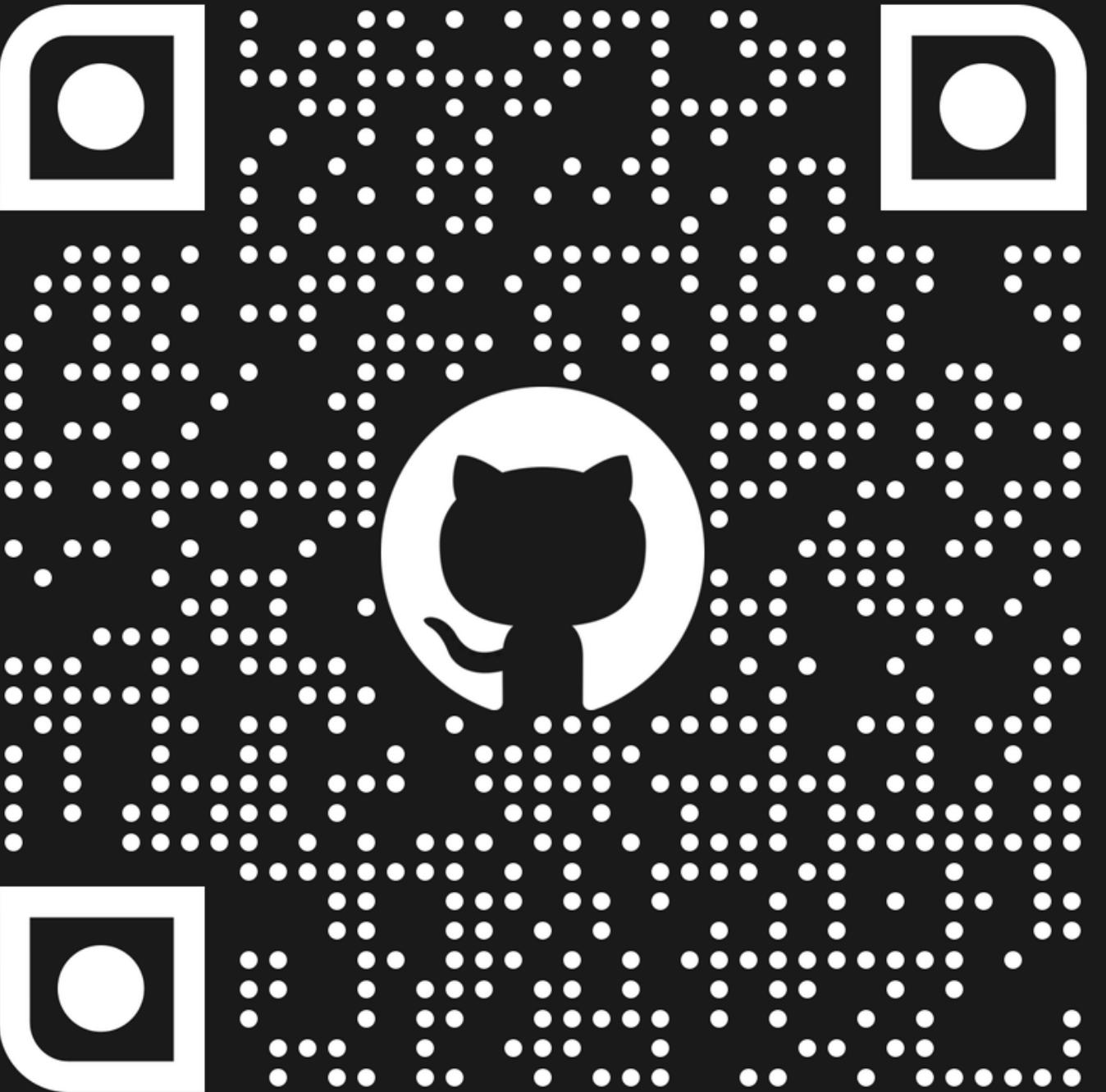
- BRIEF
- DATASET DESCRIPTION
- FEATURES ILLUSTRATION
- IMPLEMENTATION DETAILS
- RESULT DETAILS

IMAGE DATASET

- BRIEF
- DATASET DESCRIPTION
- IMPLEMENTATION DETAILS
- RESULT DETAILS



1



OR

[Click Here](#)



2

# NUMERICAL DATASET

Life Expectancy



• BRIEF

THIS REPORT SHOWCASES A  
MACHINE LEARNING MODEL  
DESIGNED TO FORECAST  
**LIFE EXPECTANCY.**  
THE PREDICTIVE MODEL  
EMPLOYS BOTH  
**LINEAR REGRESSION**  
AND  
**KNN**  
**ALGORITHMS FOR THE TASK**

# • DATASET DESCRIPTION

This study aims to address crucial gaps in past life expectancy research by considering previously overlooked factors such as immunization and human development index. Unlike prior studies that focused mainly on demographic and income-related variables, this research emphasizes a broader range of factors, spanning immunization, mortality, economy, and social aspects.

The primary objective is to leverage machine learning models, including mixed effects and multiple linear regression, to identify factors contributing to variations in life expectancy. By analyzing these diverse factors, the study aims to assist countries in pinpointing areas of focus to effectively enhance their population's life expectancy.

- **DATASET DESCRIPTION**

## Before Preprocessing

**SIZE OF DATA : 333.44 KB**

**NUMBER OF FEATURES : 22 FEATURE**

**NUMBER OF SAMPLES : 2937 SAMPLE**

# **FEATURES**

# **ILLUSTRATION**



# FEATURES ILLUSTRATION

- **Country:** The Name of The Country
- **Year:** The Year of Study
- **Status:** Developed or Developing Country status
- **Life Expectancy:** Life Expectancy in Age
- **Adult Mortality:** Adult Mortality Rates of both Sexes
- **Infant Deaths:** No. of Infant Deaths per 1000 population
- **Alcohol:** Alcohol, recorded per capita (15+) consumption (in litres of pure alcohol)
- **Percentage Expenditure:** Expenditure on health as a percentage of Gross Domestic Product per capita(%)

# FEATURES ILLUSTRATION

- **Hepatitis B:** Hepatitis B immunization coverage among 1-year-olds (%)
- **Measles:** Measles – number of reported cases per 1000 population
- **BMI:** Average Body Mass Index of entire population
- **Under-five deaths:** Number of under-five deaths per 1000 population
- **Polio:** Polio (Pol3) immunization coverage among 1-year-olds (%)

# FEATURES ILLUSTRATION

- **Total expenditure:** General government expenditure on health as a percentage of total government expenditure
- **Diphtheria:** Diphtheria tetanus toxoid and pertussis (DTP3) immunization coverage among 1-year-olds (%)
- **HIV/AIDS:** Deaths per 1 000 live births HIV/AIDS (0-4 years)
- **GDP:** Gross Domestic Product per capita (in USD)
- **Population :** Population of the country
- **Thinness 1-19 years:** Prevalence of thinness among children and adolescents for Age 10 to 19 ( % )

# FEATURES ILLUSTRATION

- **thinness 5–9 years:** Prevalence of thinness among children for Age 5 to 9(%)
- **Income composition of resources:** Human Development Index in terms of income composition of resources (index ranging from 0 to 1)
- **Schooling :** Number of years of Schooling(years)

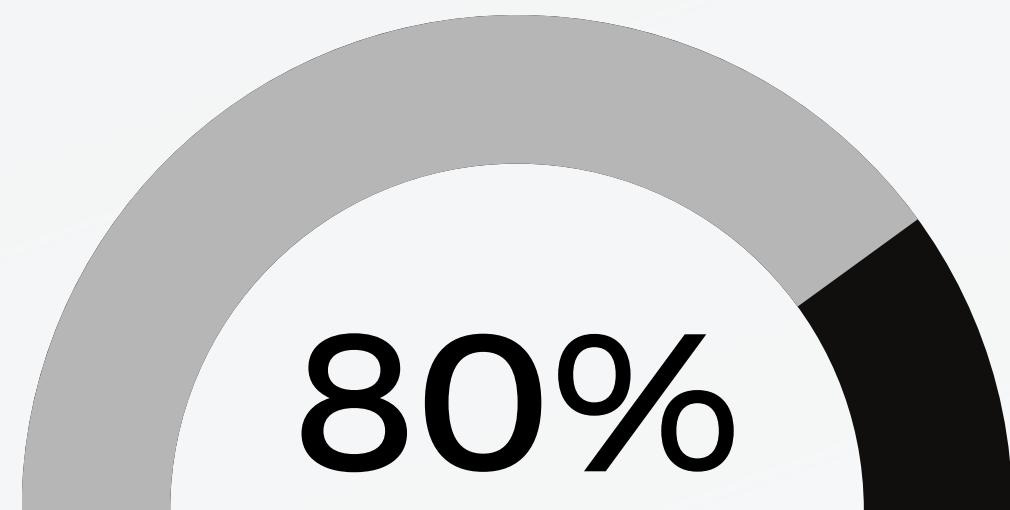
## DATA TYPES

**float64(16), int64(4), object(2)**

# **IMPLEMENTATION DETAILS**



# DATA SPLITTING



TRAIN



TEST

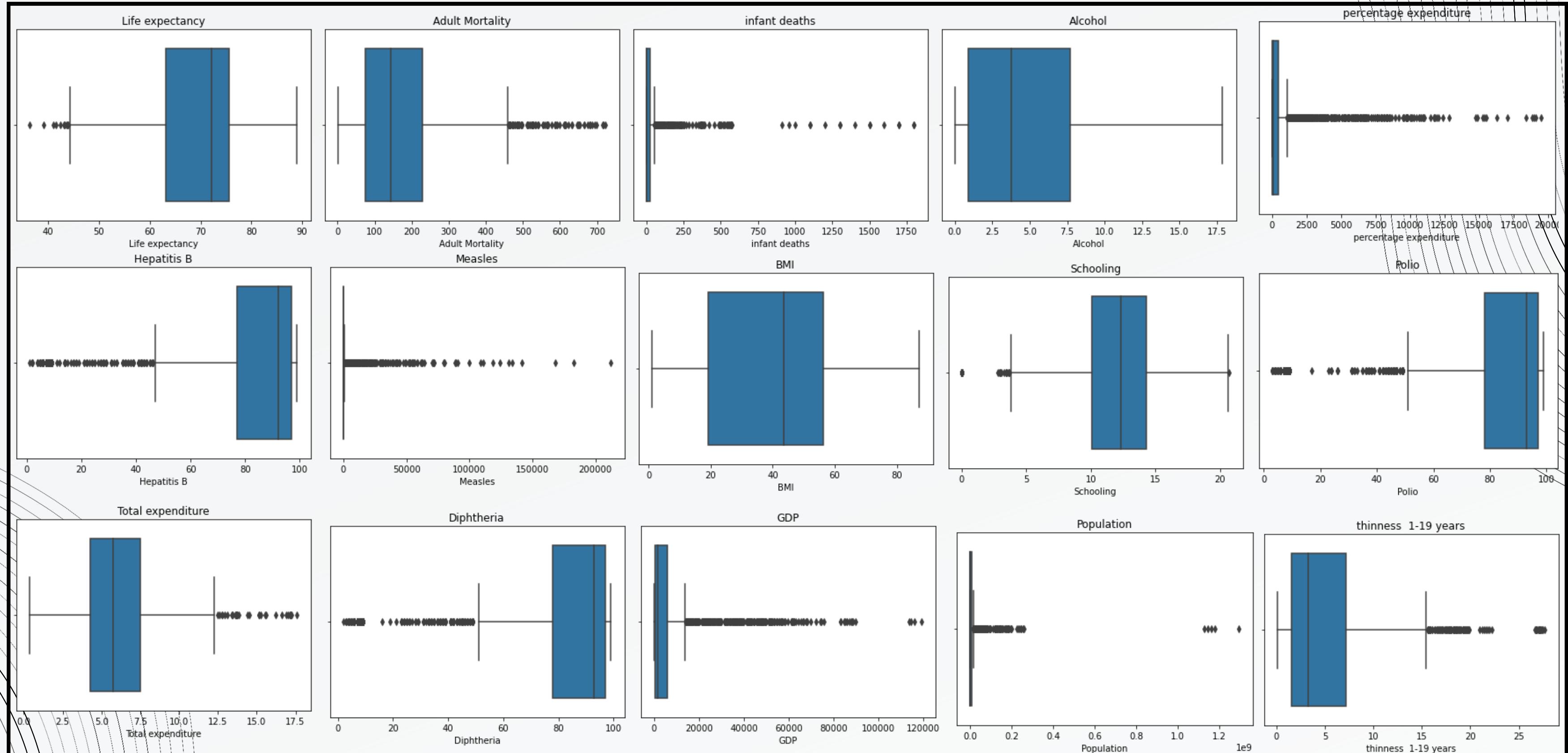
# ANALYZING

- Data Import and Preprocessing:
  - Reads and explores the dataset using Pandas, including examining data samples, information, and statistical summaries.
  - Handles null values, strips column names, and visualizes distributions and outliers in numerical features.
  - Explores categorical features with bar and pie charts, depicting counts and proportions.
- Exploratory Data Analysis (EDA):
  - Investigates relationships between features and the target variable (life expectancy).
  - Visualizes temporal trends in life expectancy and infant deaths over the years.
  - Studies correlations among numerical features and their impact on life expectancy.

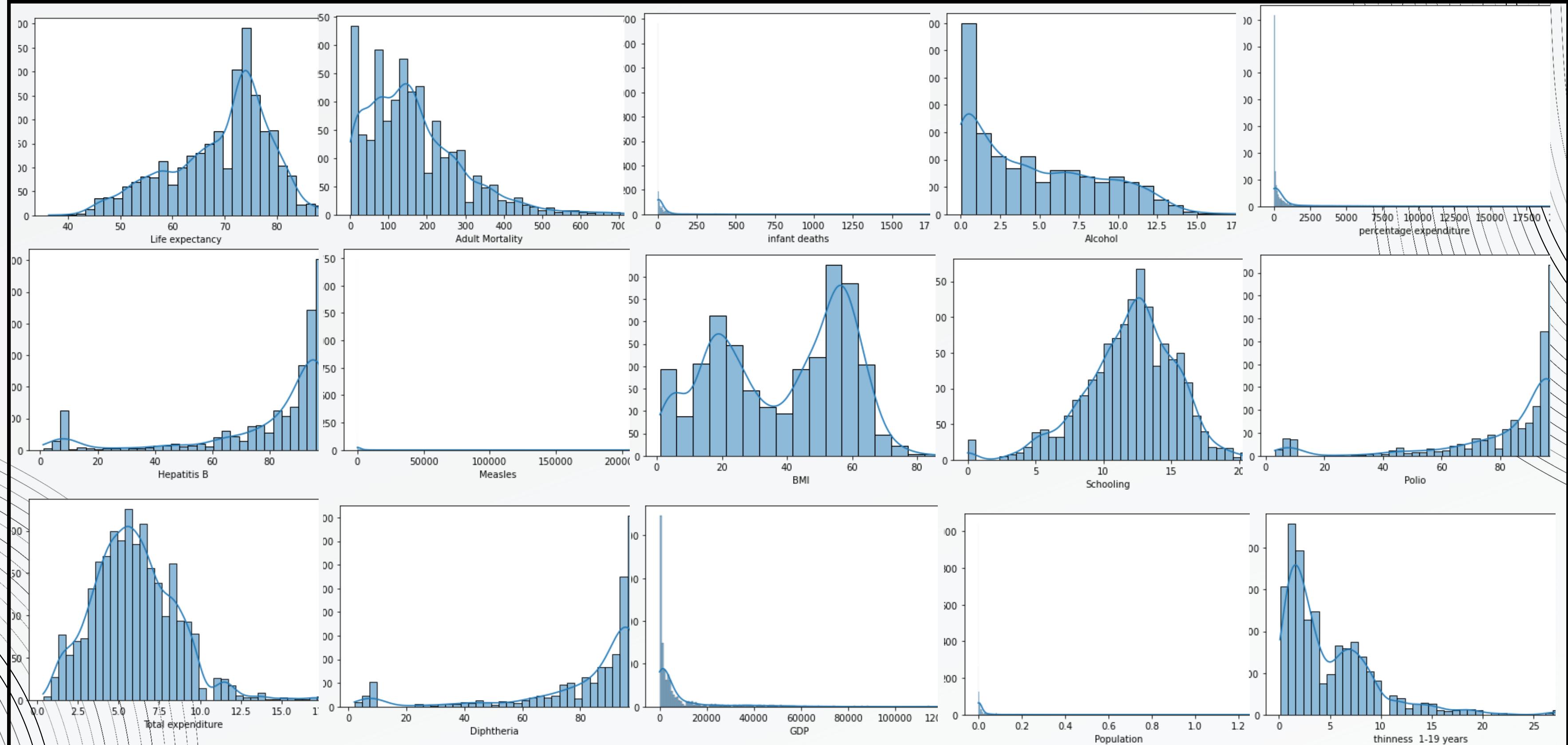
# ANALYZING

- Data Preprocessing and Feature Engineering:
  - Addresses null values by imputing with median values.
  - Handles outliers using an Interquartile Range (IQR) method.
  - Drops specific columns and encodes categorical features for modeling.
- Modeling:
  - Splits data into training and testing sets.
  - Standardizes numerical features.
  - Applies Linear Regression for predicting life expectancy, evaluates model performance, and visualizes predictions against actual values.
  - Determines feature importance using coefficient analysis and visualizes the linear regression feature importance.
- Model Evaluation and Improvement:
  - Assesses model performance using Root Mean Squared Error (RMSE), R-squared score, and cross-validation.
  - Implements K-Nearest Neighbors (KNN) Regression, optimizing for the best 'k' value using GridSearchCV, and evaluates its performance.

# OUTLIERS Before...



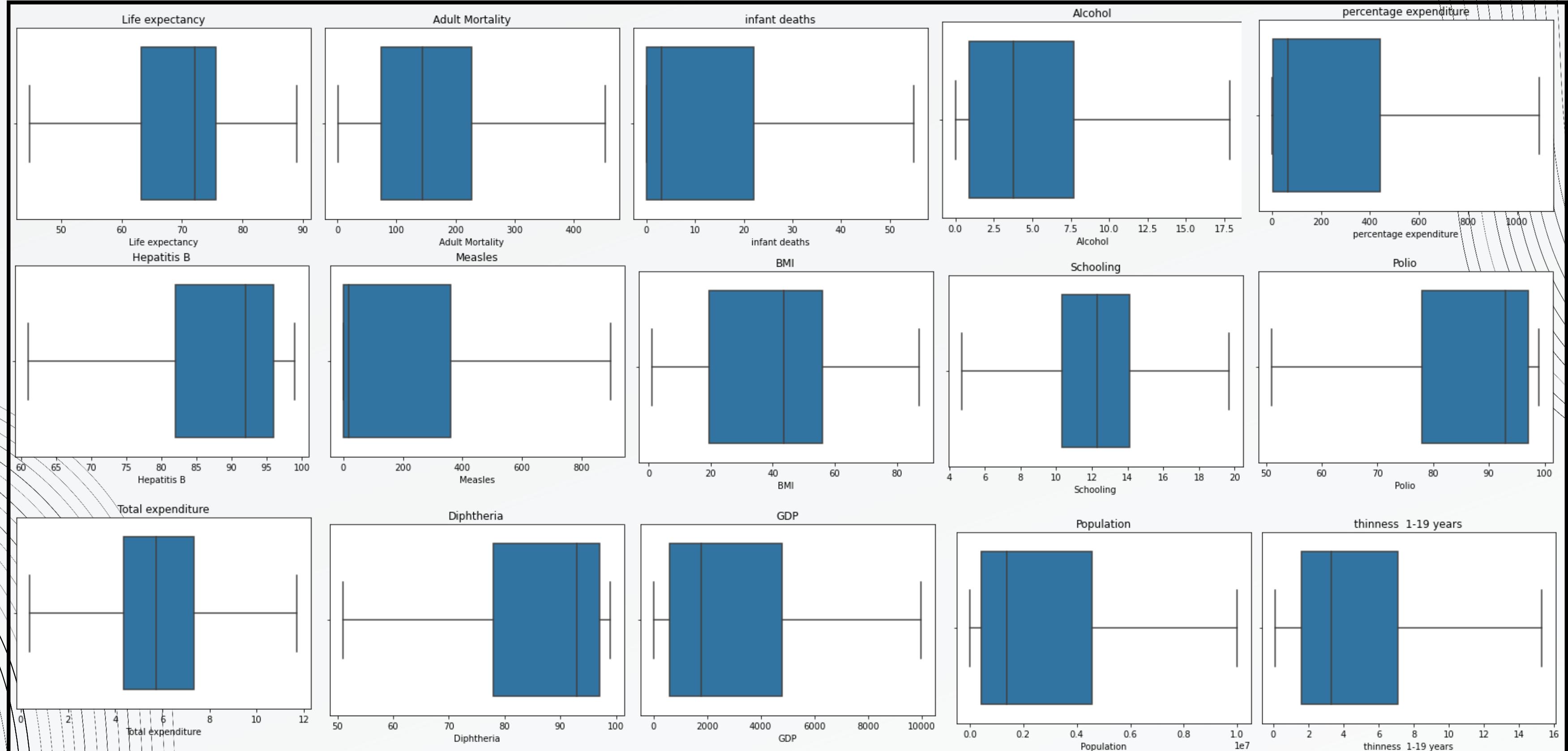
# DISTRIBUTION Before...



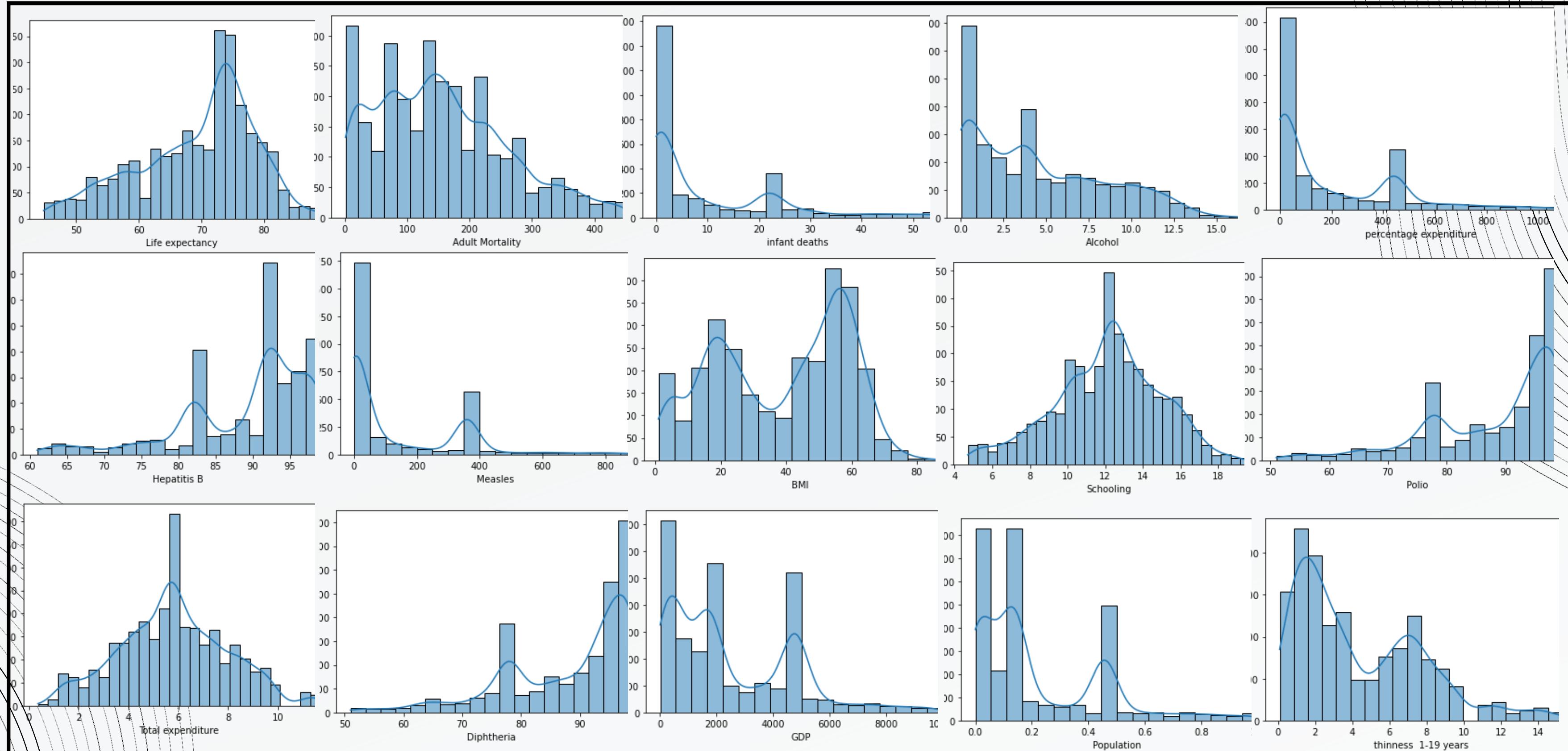
# RESULT DETAILS



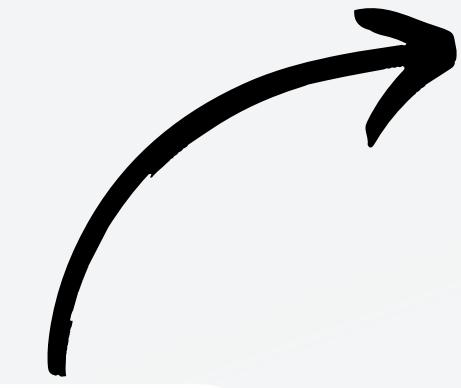
# OUTLIERS After...



# DISTRIBUTION After...



## LINEAR REGRESSION



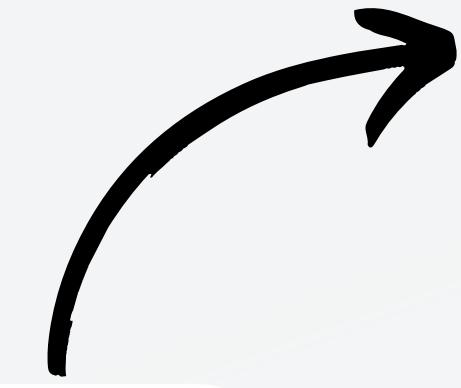
## CROSS VALIDATION



- Train accuracy 79%
- Test accuracy 80%
- RMSE for training : 4.23
- RMSE for testing : 4.16
- R2\_score 88%

- Cross Validation score : 78.4%
- The training/validation ratio is 4:1 for each fold

KNN



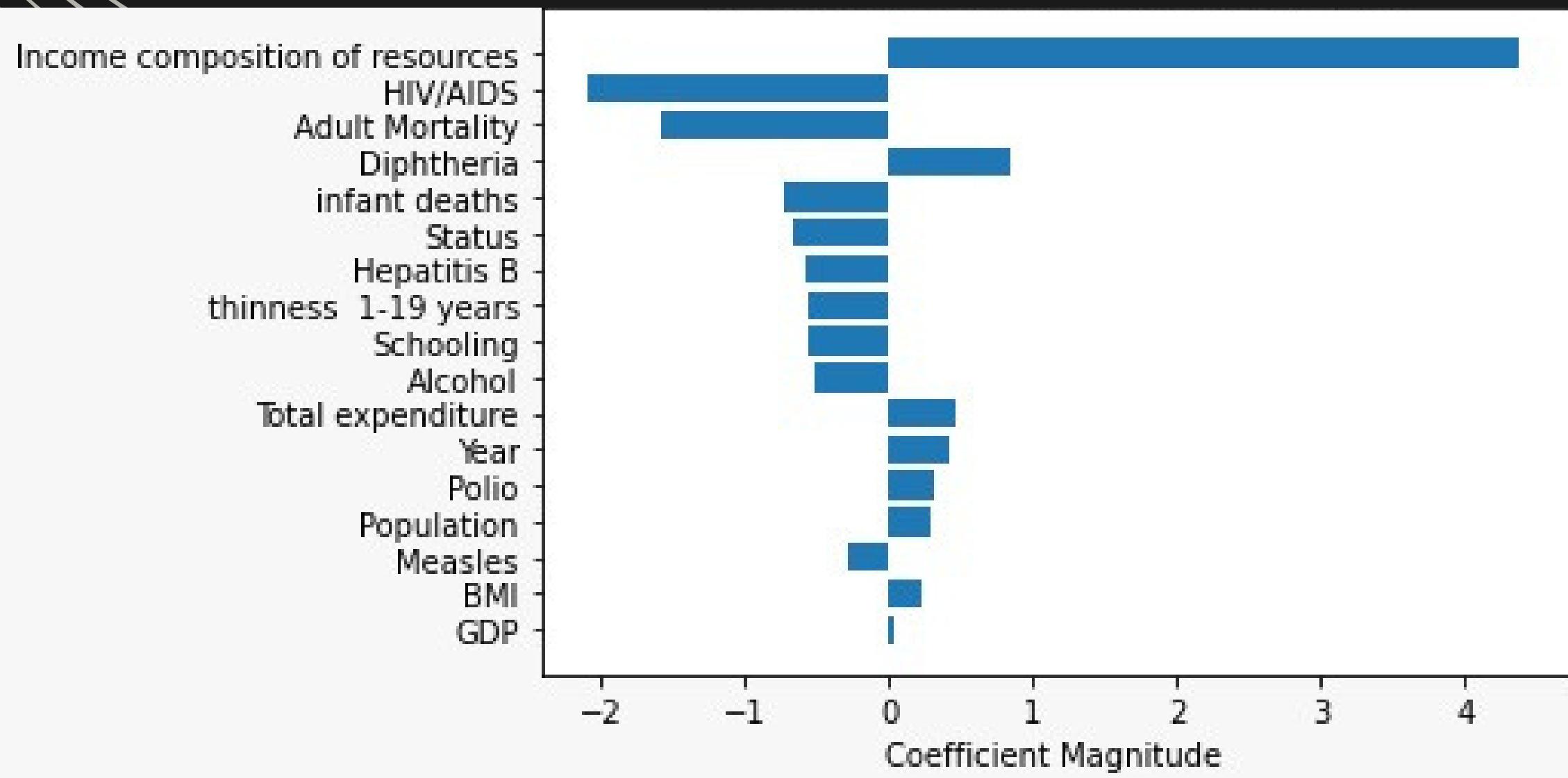
- Train accuracy  94%
- Test accuracy  88%
- RMSE for training : 2.21
- RMSE for testing : 3.23
- n\_neighbors : 3

CROSS VALIDATION



- Cross Validation score : 86%
- The training/validation ratio is 4:1 for each fold

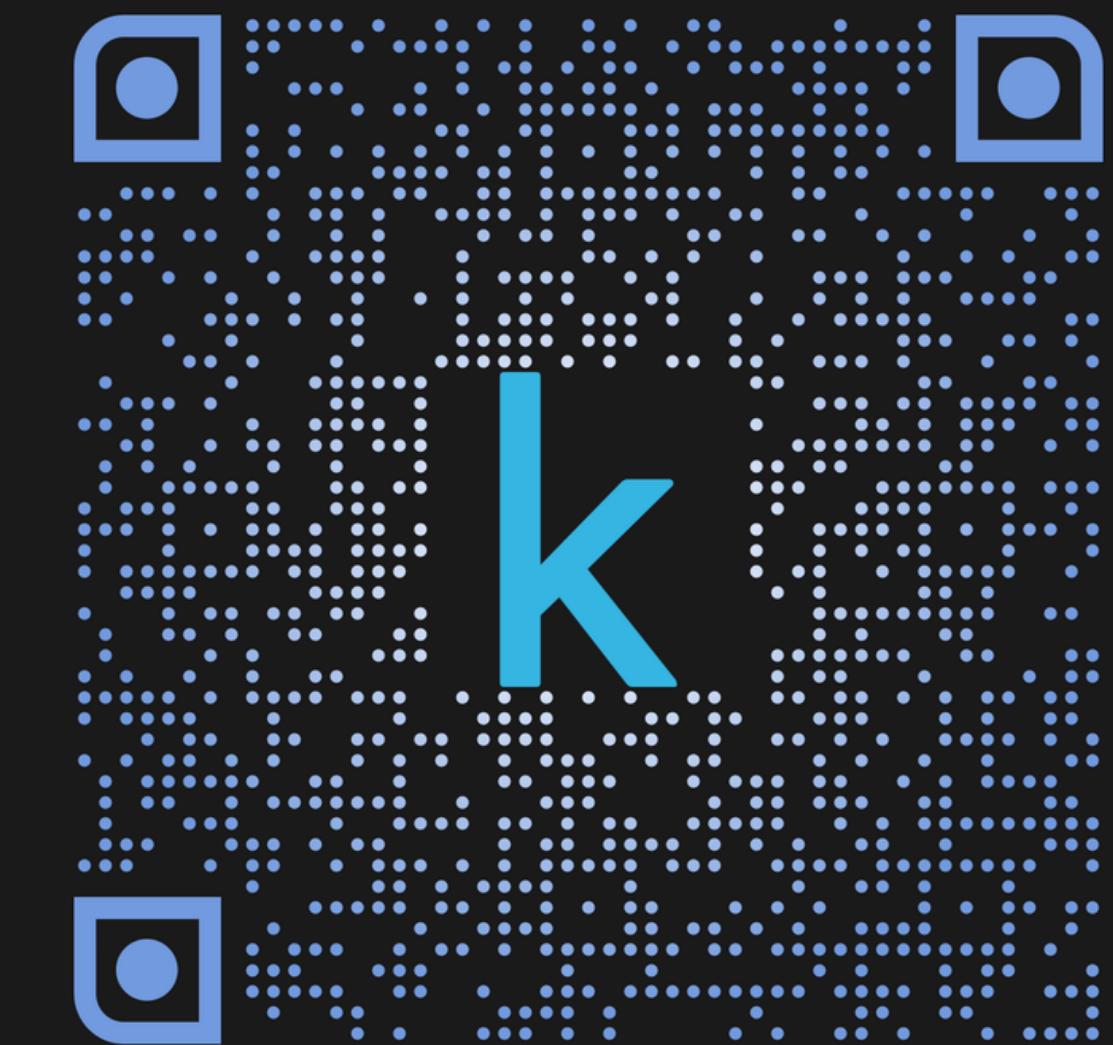
# FEATURES IMPORTANCE



3

# IMAGE DATASET

Stanford Dogs Dataset

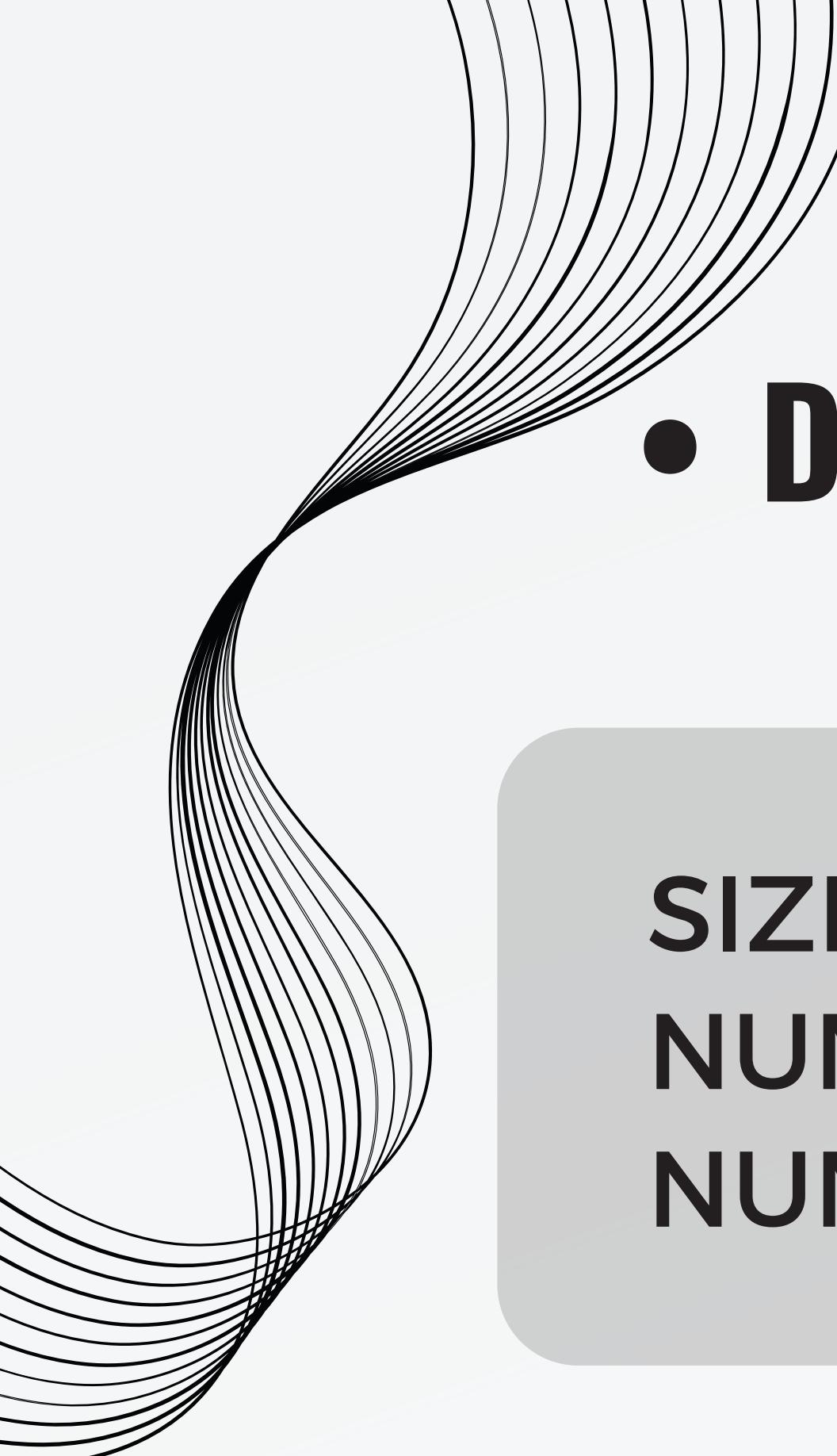


• BRIEF

THIS SECTION SHOWCASES A  
MACHINE LEARNING MODEL  
DESIGNED TO FORECAST  
**DIVERSE IMAGES OF DOGS  
ACROSS 120 BREEDS.**  
THE PREDICTIVE MODEL  
EMPLOYS BOTH  
**LOGISTIC REGRESSION**  
AND  
**KMEANS**  
**ALGORITHMS FOR THE TASK**

## • DATASET DESCRIPTION

The Stanford Dogs dataset contains images of 120 breeds of dogs from around the world. This dataset has been built using images and annotation from ImageNet for the task of fine-grained image categorization. It was originally collected for fine-grain image categorization, a challenging problem as certain dog breeds have near identical features or differ in colour and age.



## • DATASET DESCRIPTION

### Before Preprocessing

**SIZE OF DATA : 788.05 MB**

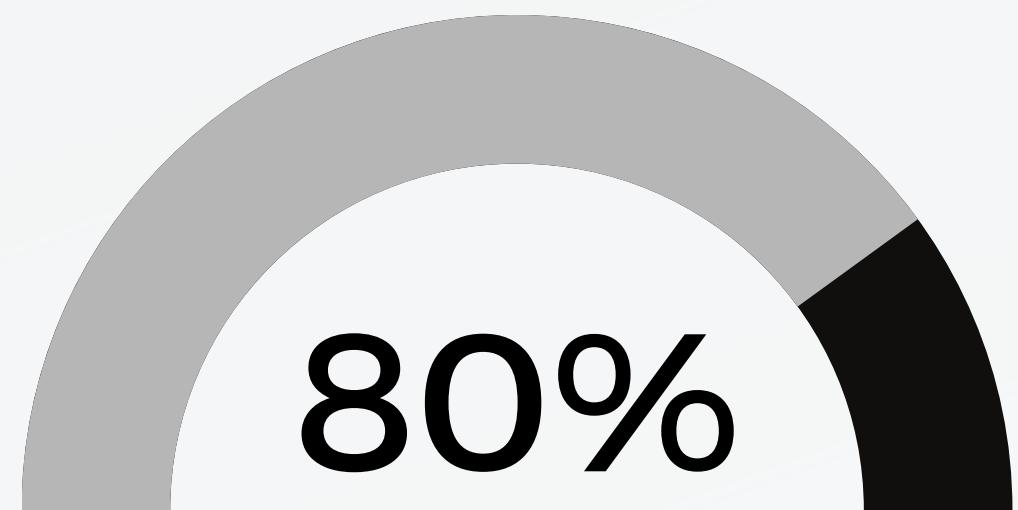
**NUMBER OF CLASSES: 33 CLASS**

**NUMBER OF IMAGES: 20,580**

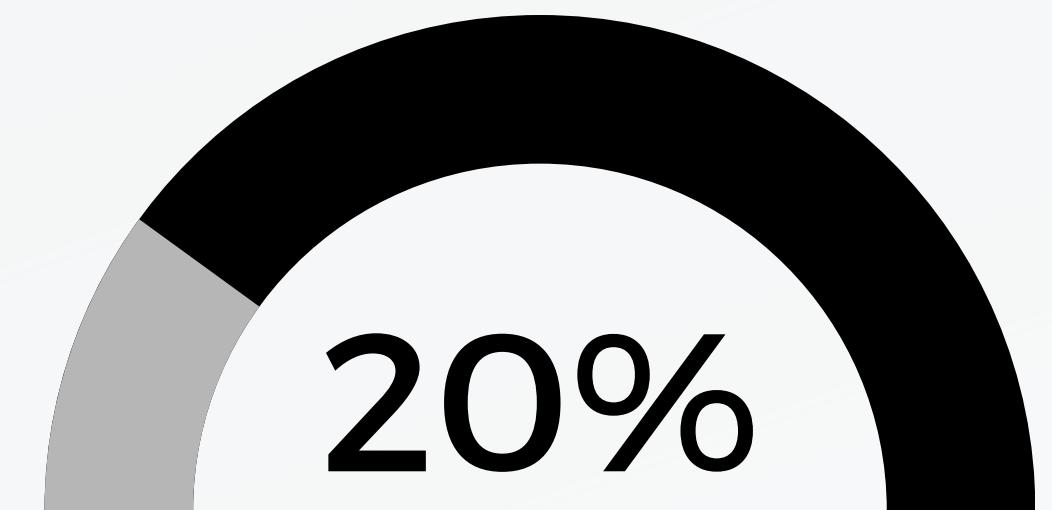
# **IMPLEMENTATION DETAILS**



# DATA SPLITTING



TRAIN



TEST

# ANALYZING

1. **Data Preparation:** The code reads images and annotations from specified directories, extracts bounding box coordinates, and preprocesses the images for analysis.
2. **K-Means Clustering:** Utilizes K-Means clustering to cluster flattened image arrays based on pixel values, generating clusters and visualizes them using Principal Component Analysis (PCA).
3. **SIFT Feature Extraction:** Extracts SIFT (Scale-Invariant Feature Transform) features from the images, displays original and resized images with detected SIFT features, and performs clustering based on SIFT descriptors.

# ANALYZING

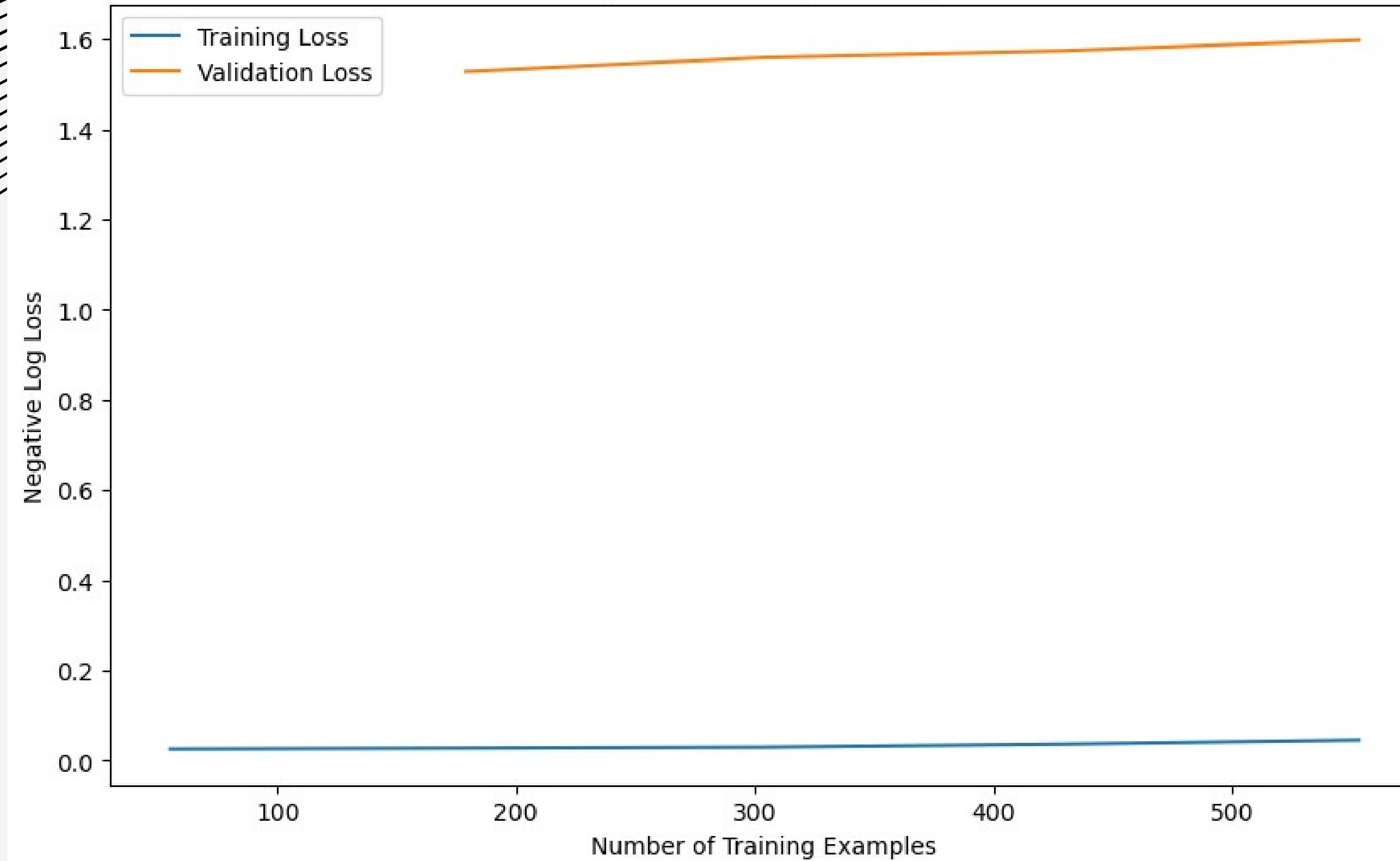
1. **SIFT Features Analysis:** Processes the SIFT descriptors, pads them for uniformity, performs K-Means clustering on the modified descriptors, and displays images from different clusters to visualize grouping.
2. **Clustering Evaluation:** Evaluates clusters' characteristics, including the number of images in each cluster, using indices and visualizations.
3. **Elbow Method and PCA:** Employs the Elbow Method to determine an optimal number of clusters, then applies PCA to visualize the clusters in a 2D space.

# RESULT DETAILS

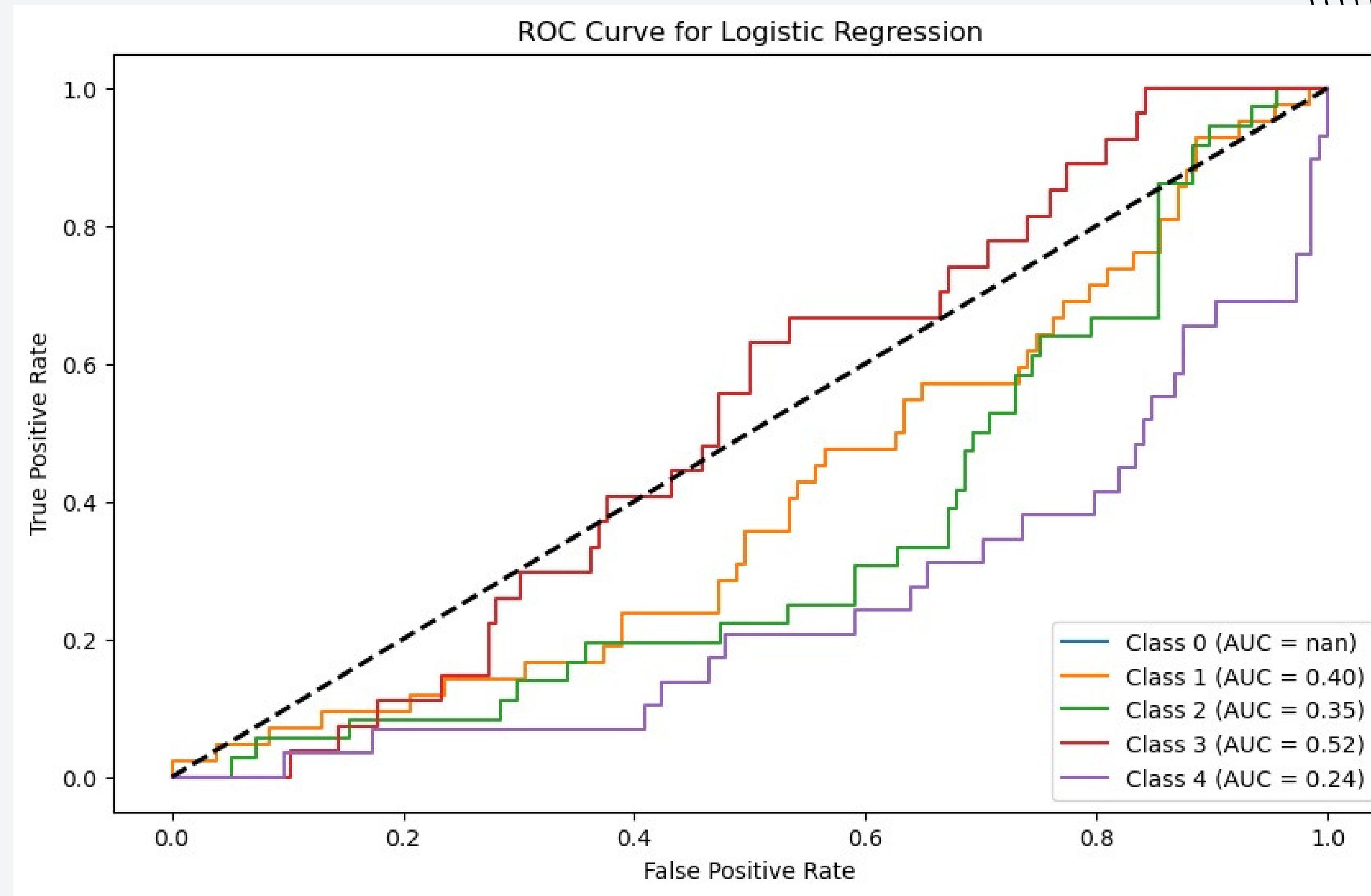


# LOSS CURVE

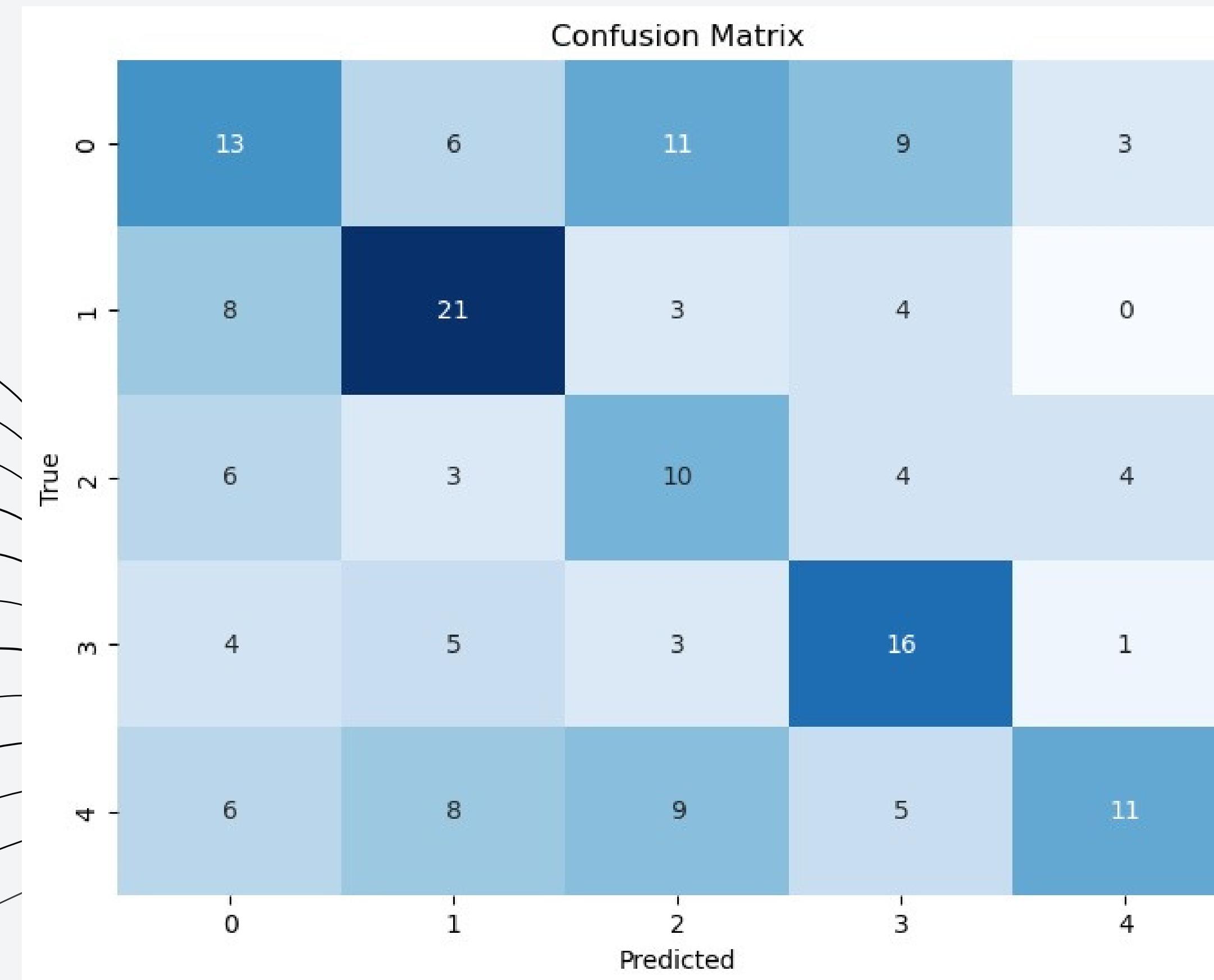
Learning Curve for Logistic Regression



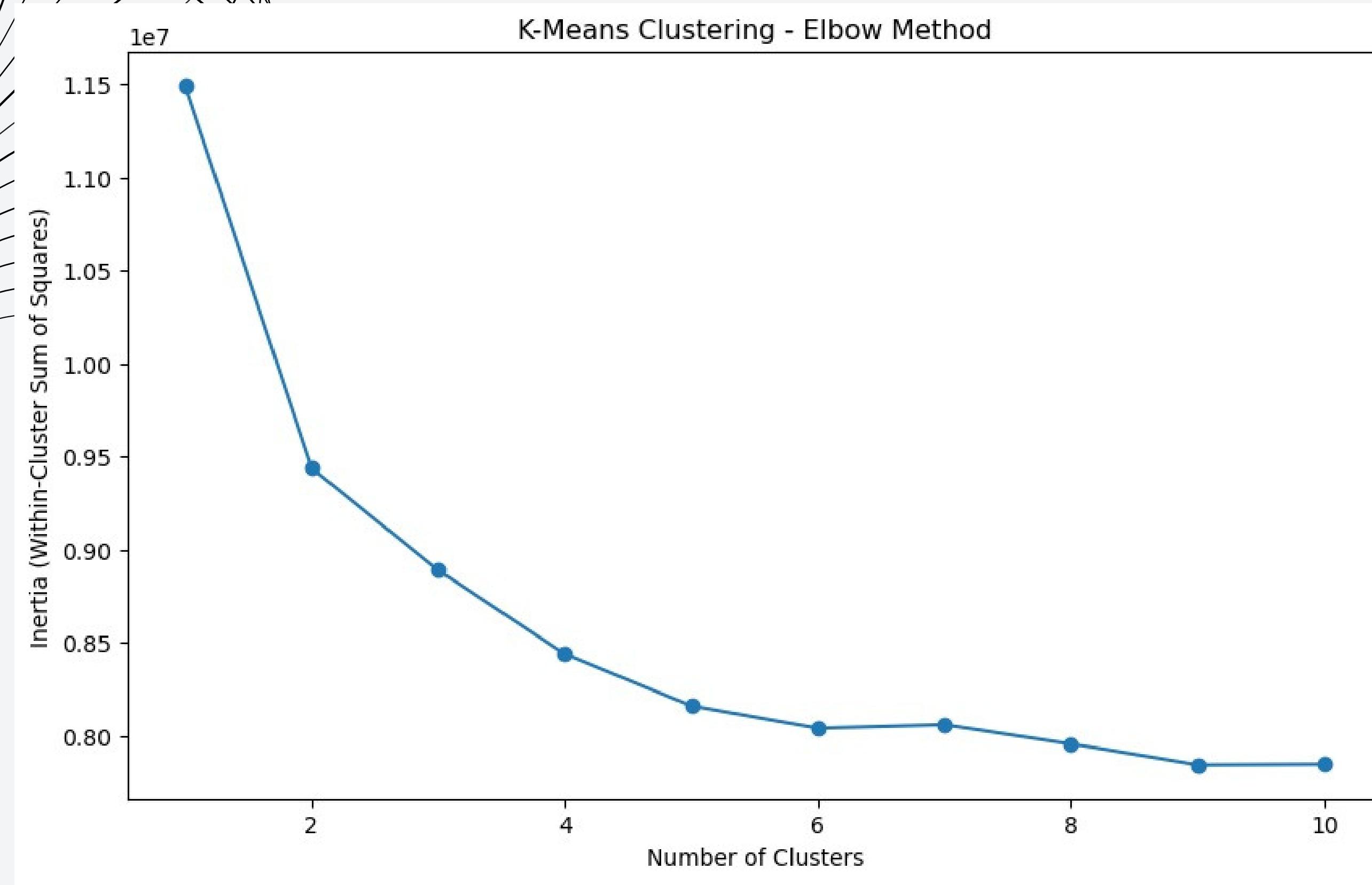
# ROC CURVE



# CONFUSSION MATRIX

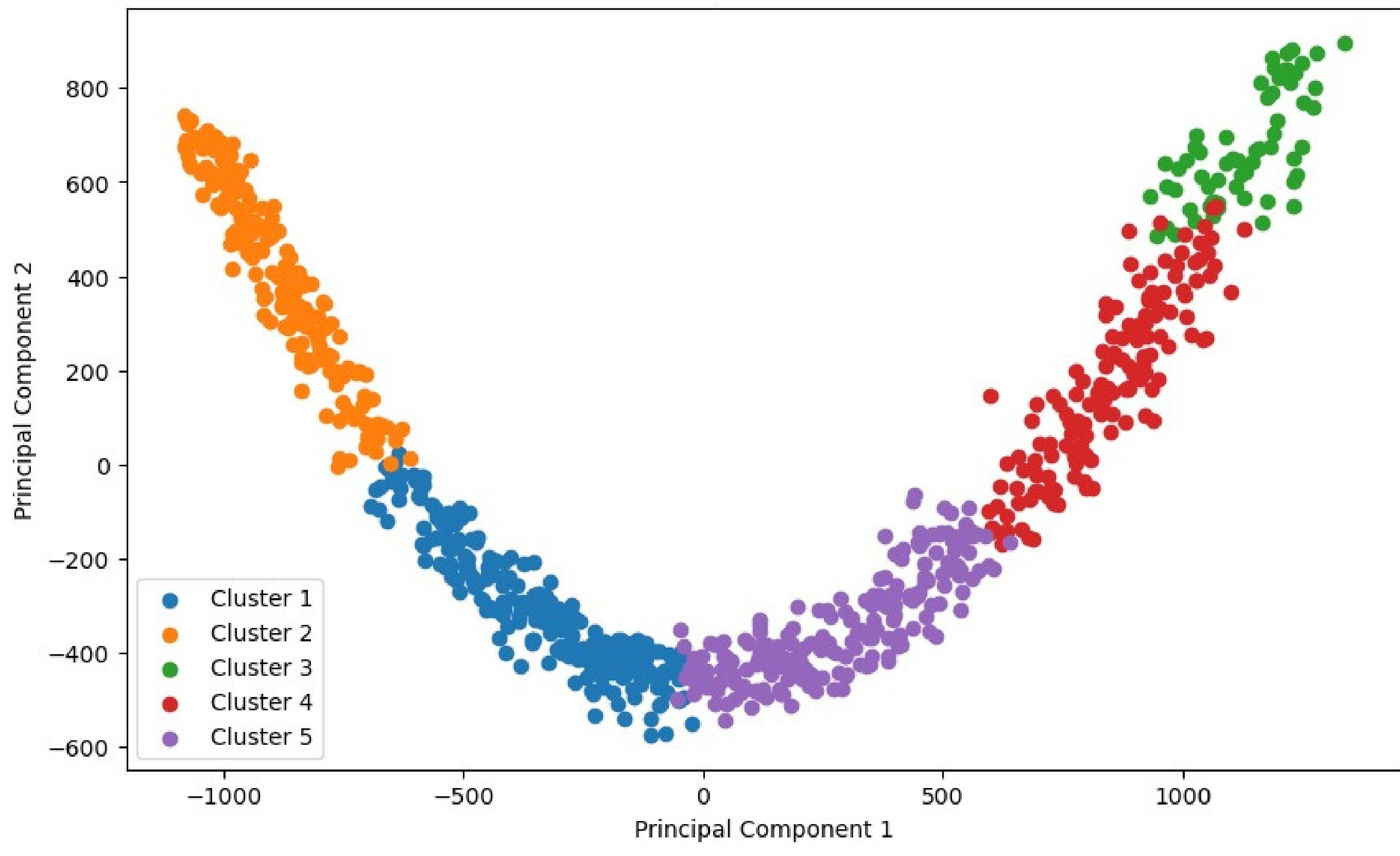


# ELBOW METHOD



# SCATTER PLOT

Scatterplot of Clusters



# AFTER CLUSTERING

Cluster 1, Image 605



Cluster 1, Image 829



Cluster 1, Image 129



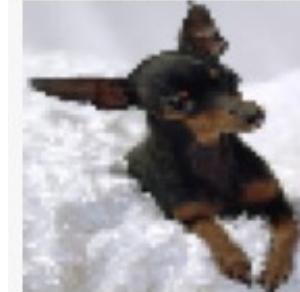
Cluster 1, Image 627



Cluster 1, Image 67



Cluster 2, Image 752



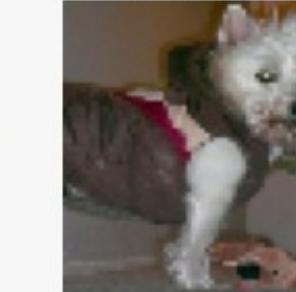
Cluster 2, Image 414



Cluster 2, Image 95



Cluster 2, Image 287



Cluster 2, Image 265



Cluster 3, Image 656



Cluster 3, Image 726



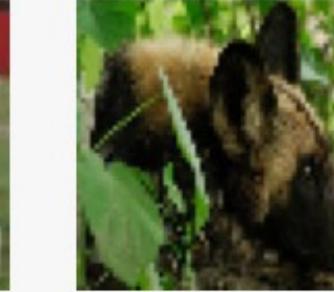
Cluster 3, Image 785



Cluster 3, Image 138



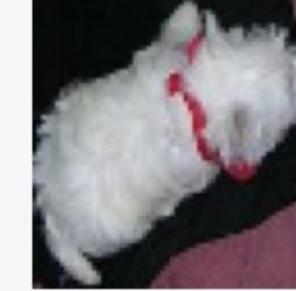
Cluster 3, Image 709



Cluster 4, Image 892



Cluster 4, Image 370



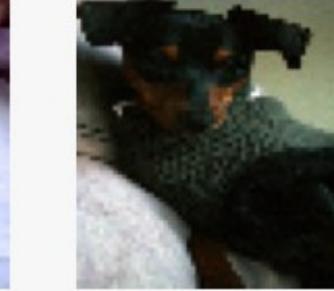
Cluster 4, Image 399



Cluster 4, Image 779



Cluster 4, Image 773



Cluster 5, Image 22



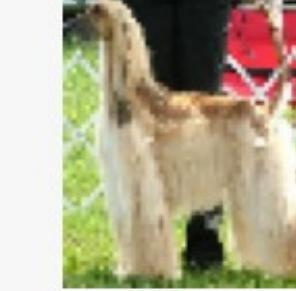
Cluster 5, Image 824



Cluster 5, Image 126



Cluster 5, Image 152

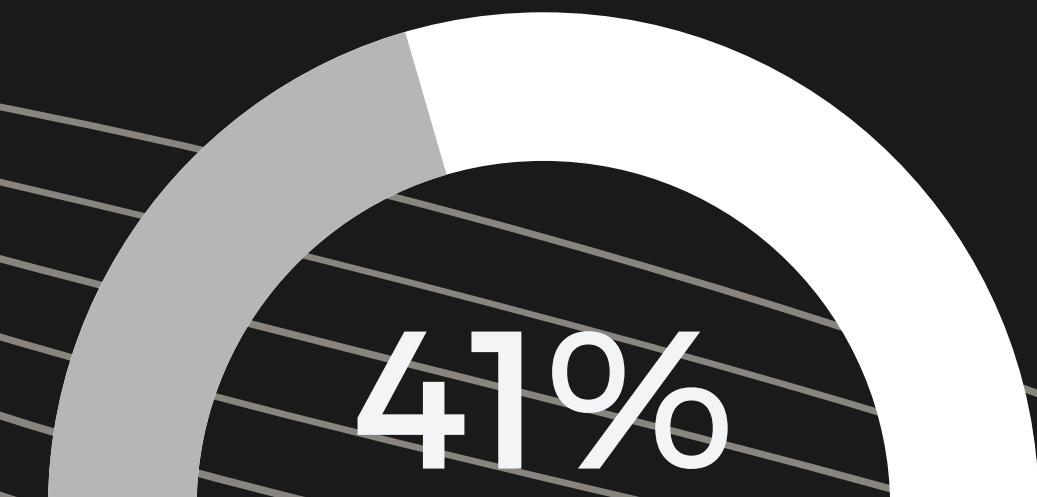


Cluster 5, Image 889



# ACCURACY

LOGISTIC REGRESSION



41%

A semi-circular gauge chart is positioned in the center-right area of the slide. It consists of a white outer ring and a dark grey inner segment. The number "41%" is displayed in white text within the dark segment. The background of the slide features a series of light grey curved lines that radiate from the bottom left towards the top right, creating a dynamic visual effect.

# OUR TEAM



**MUHAMMAD  
YASSER**



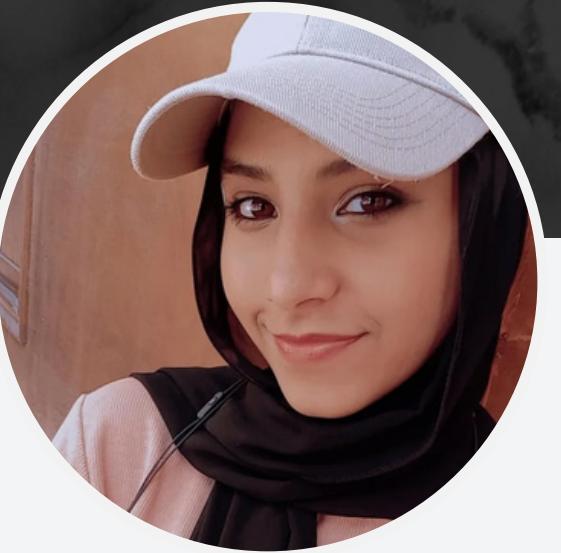
**MOHAMMAD  
TAREK**



**MOHAMED  
HUSSIAN**



**ASSER  
HASSAN**



**MADIHA  
SAEID**



**ESRAA  
MOHAMED**



**HANIA  
RUBY**

# THANKS

