



# HEART FAILURE PREDICTION USING PYSPARK

Pham Thuy Dung - 20521214

---

# AGENDA

- Introduction
- Heart Failure Prediction
- Process Flow Diagram
- PySpark Overview
- Machine Learning
- Data analysis and visualization
- Results



---

# INTRODUCTION

Heart failure prediction is an essential area of medical research and healthcare, aiming to identify individuals who are at risk of developing heart failure before they experience any symptoms.



# HEART FAILURE

- a medical condition related to the cardiovascular system
- no longer functions efficiently to pump enough blood to meet the body's needs



# HEART FAILURE PREDICTION

- process of using various medical, clinical, and lifestyle data to assess an individual's risk of developing heart failure



Loading Dataset



Data Pre-Processing



Data Transformation



Classifiers

- Logistic Regression
- Support Vector Machines
- Random Forest Classifier
- Multilayer Perceptron Classifier
- Decision Tree



Evaluating Model

## PROCESS FLOW DIAGRAM

---

# OVERVIEW OF PYSPARK

Spark SQL  
and  
DataFrames

Pandas API  
on Spark

Structured  
Streaming

Machine  
Learning  
(MLlib)

Spark Core and RDDs



# DATASET

Features	Description
<b>Age</b>	age
<b>Sex</b>	Sex(M, F)
<b>ChestPainType</b>	Chest pain type (ASY, NAP, Other)
<b>RestingBP</b>	Resting blood pressure
<b>Cholesterol</b>	Serum cholestrol
<b>FastingBS</b>	Fasting blood sugar
<b>RestingECG</b>	Resting electrocardiogram results (Normal, LVH, Others)
<b>MaxHR</b>	Maximum heart rate achieved
<b>ExerciseAngina</b>	Exercise induced angina
<b>Oldpeak</b>	Oldpeak = ST
<b>ST_Slope</b>	The slope of the peak exercise ST segment
<b>HeartDisease</b>	Target to predict

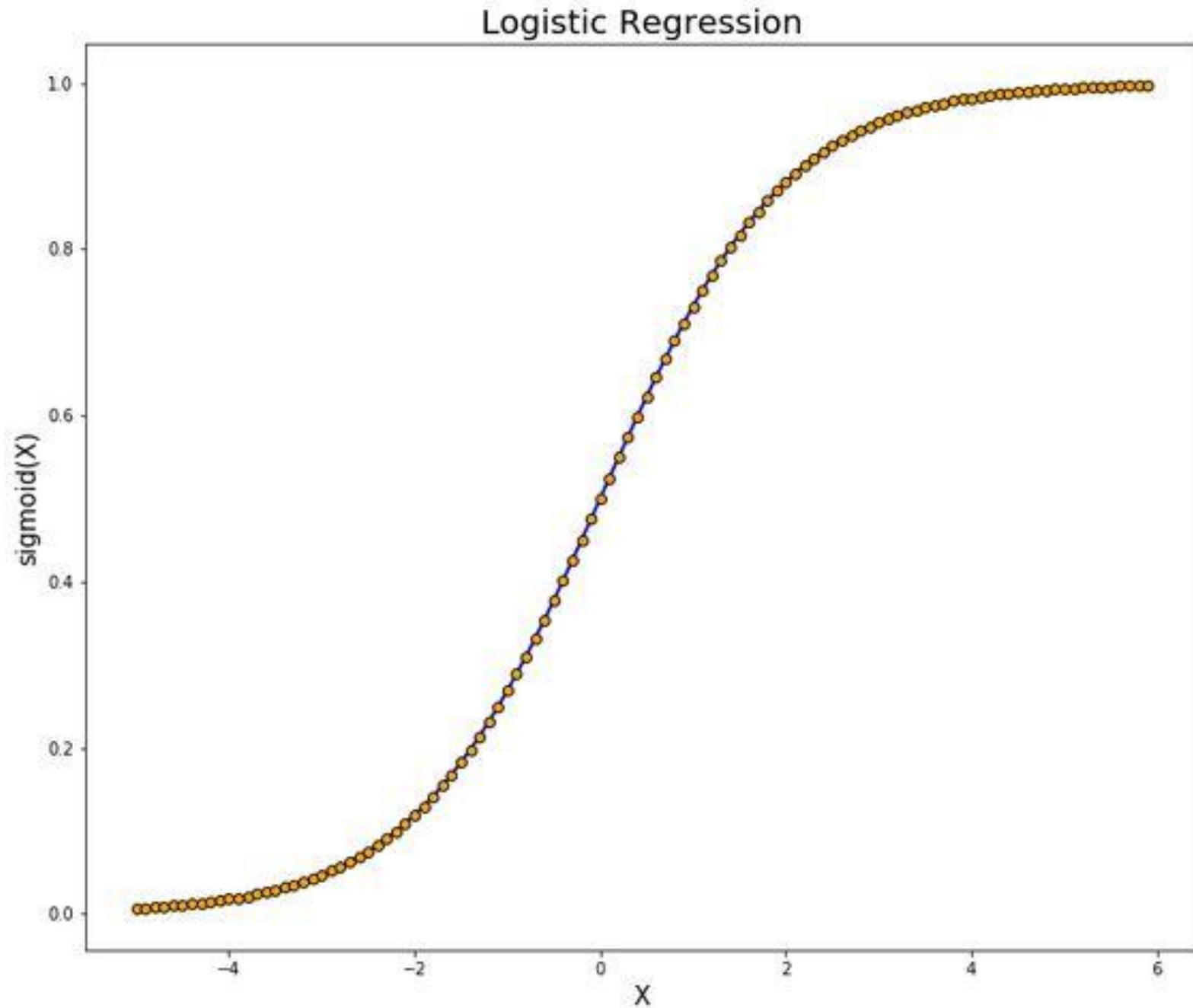


# MACHINE LEARNING

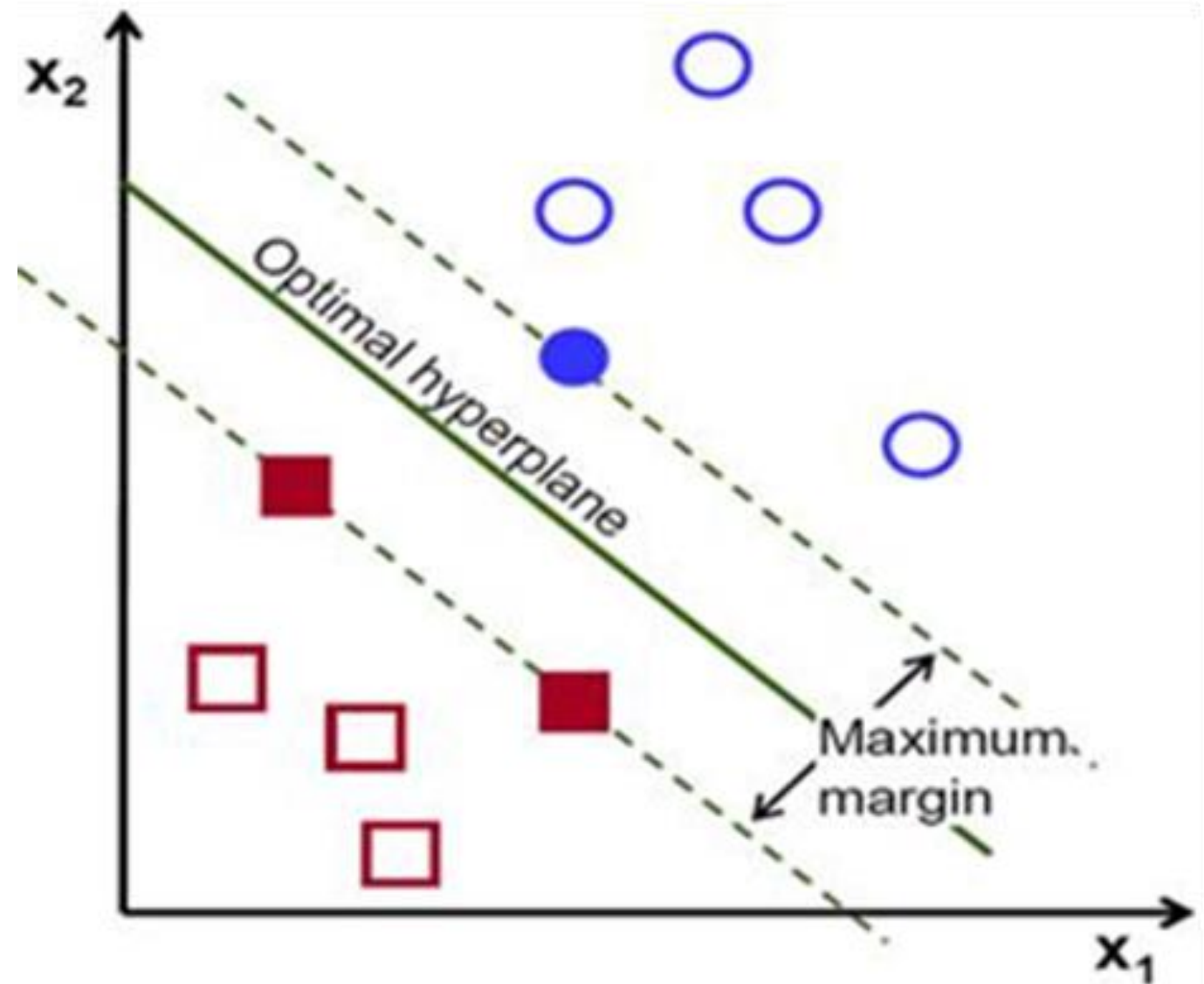


# LOGISTICS REGRESSION

$$p(x) = \frac{1}{1 + e^{-(x-\mu)/s}}$$

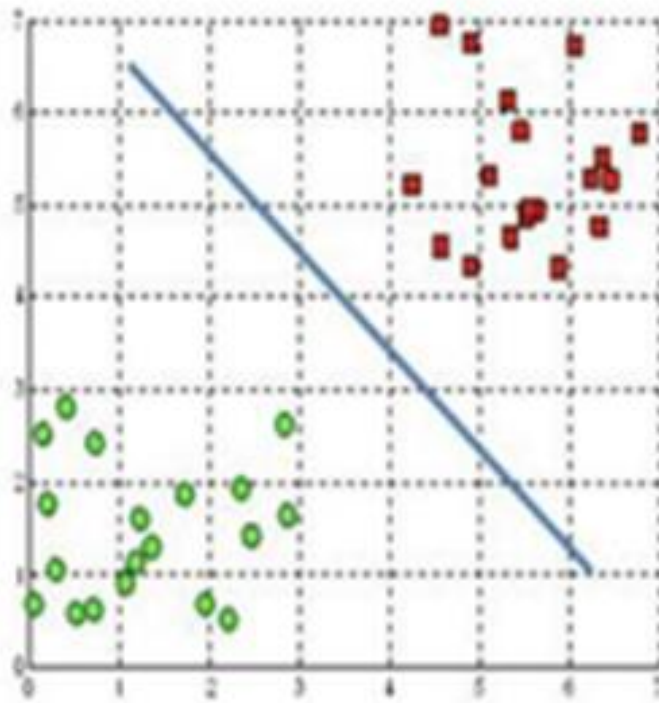


# SUPPORT VECTOR MACHINE (SVM)

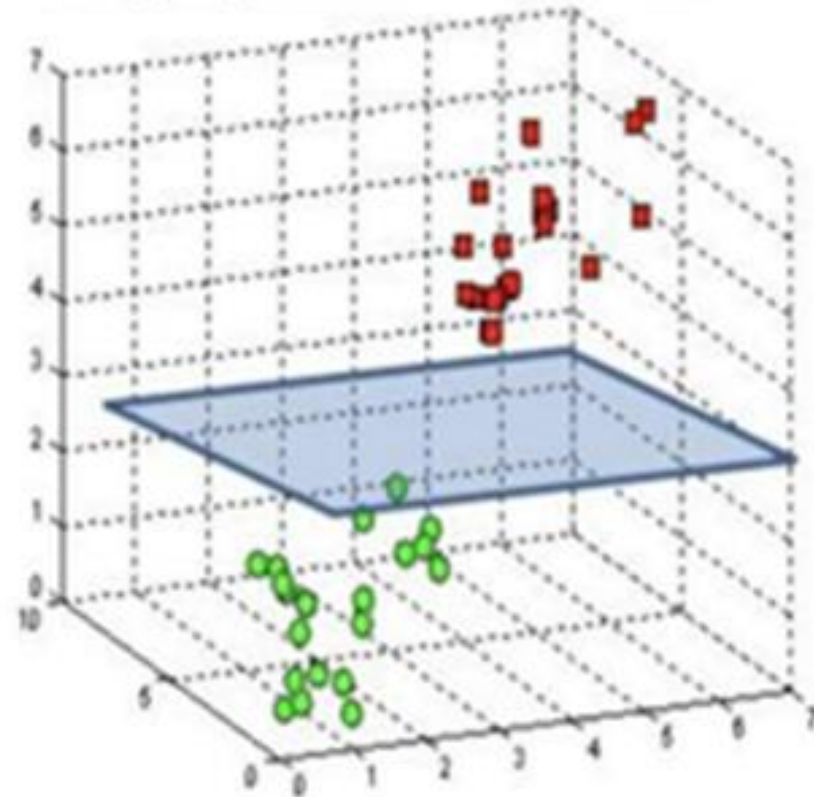


# SUPPORT VECTOR MACHINE (SVM)

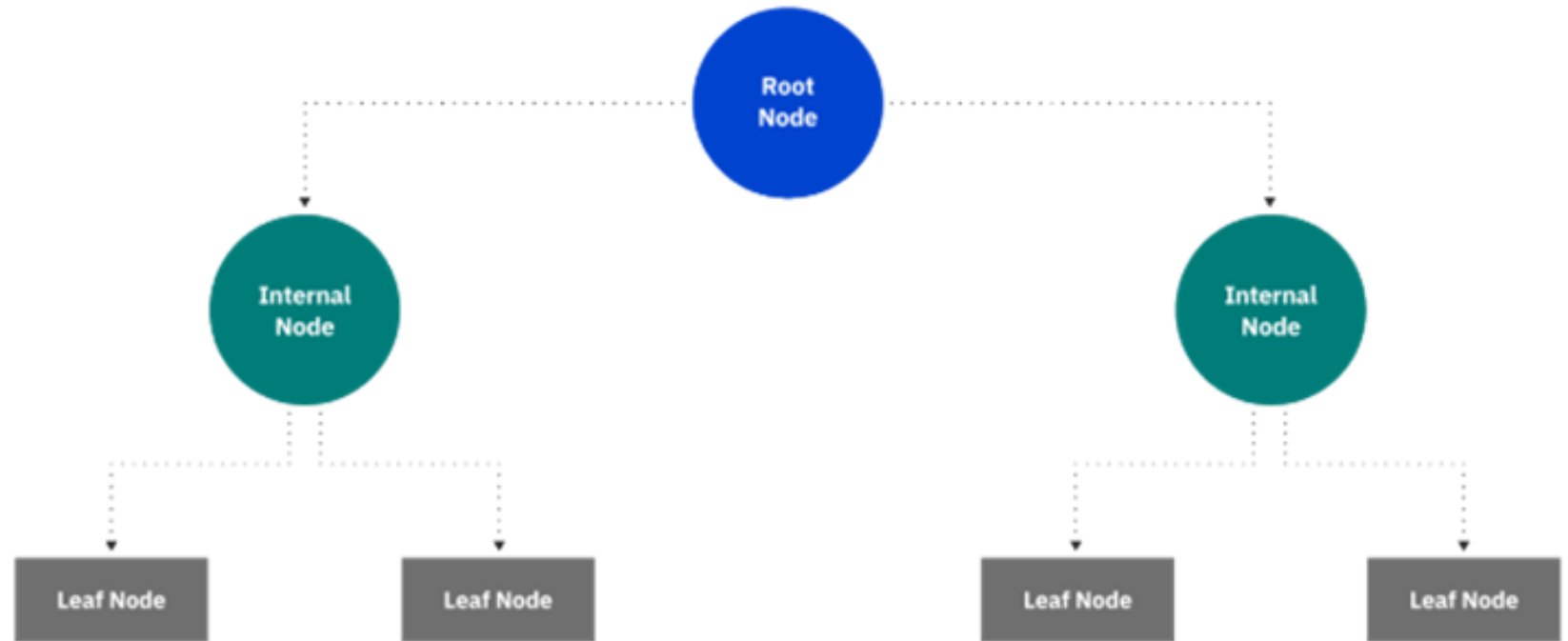
A hyperplane in  $\mathbb{R}^2$  is a line



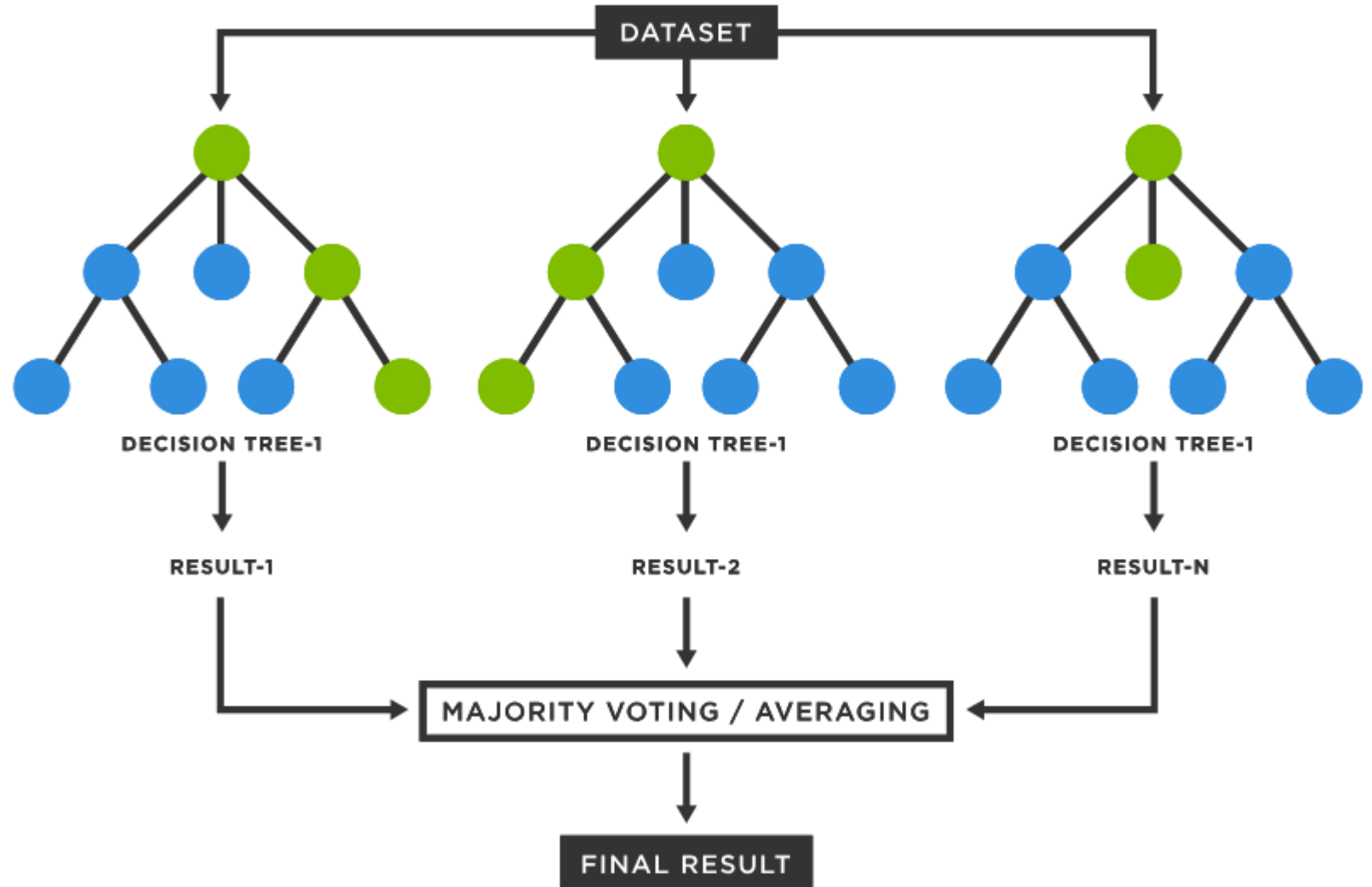
A hyperplane in  $\mathbb{R}^3$  is a plane



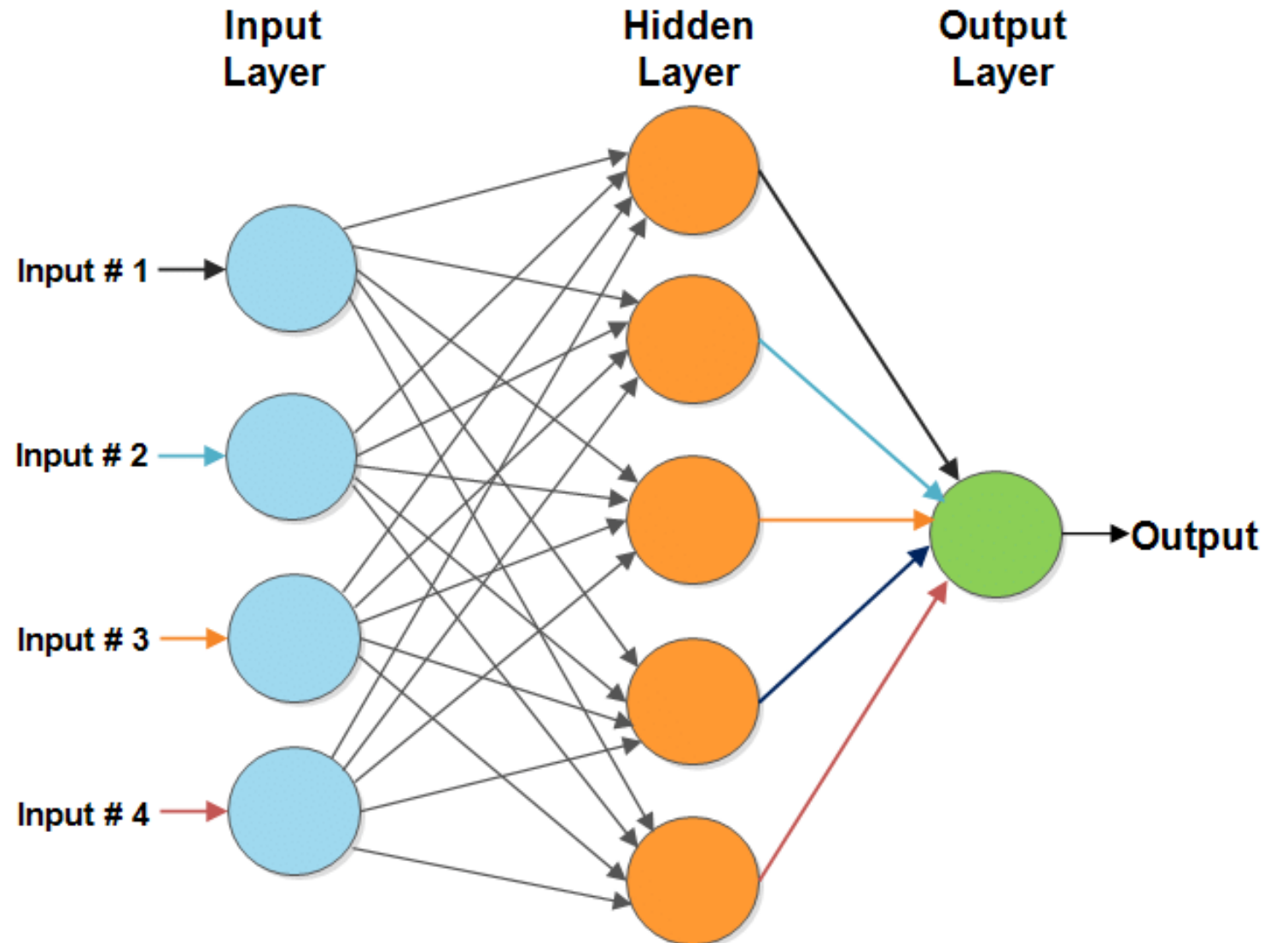
# DECISION TREE



# RANDOM FOREST



# MULTI-LAYER PERCEPTION

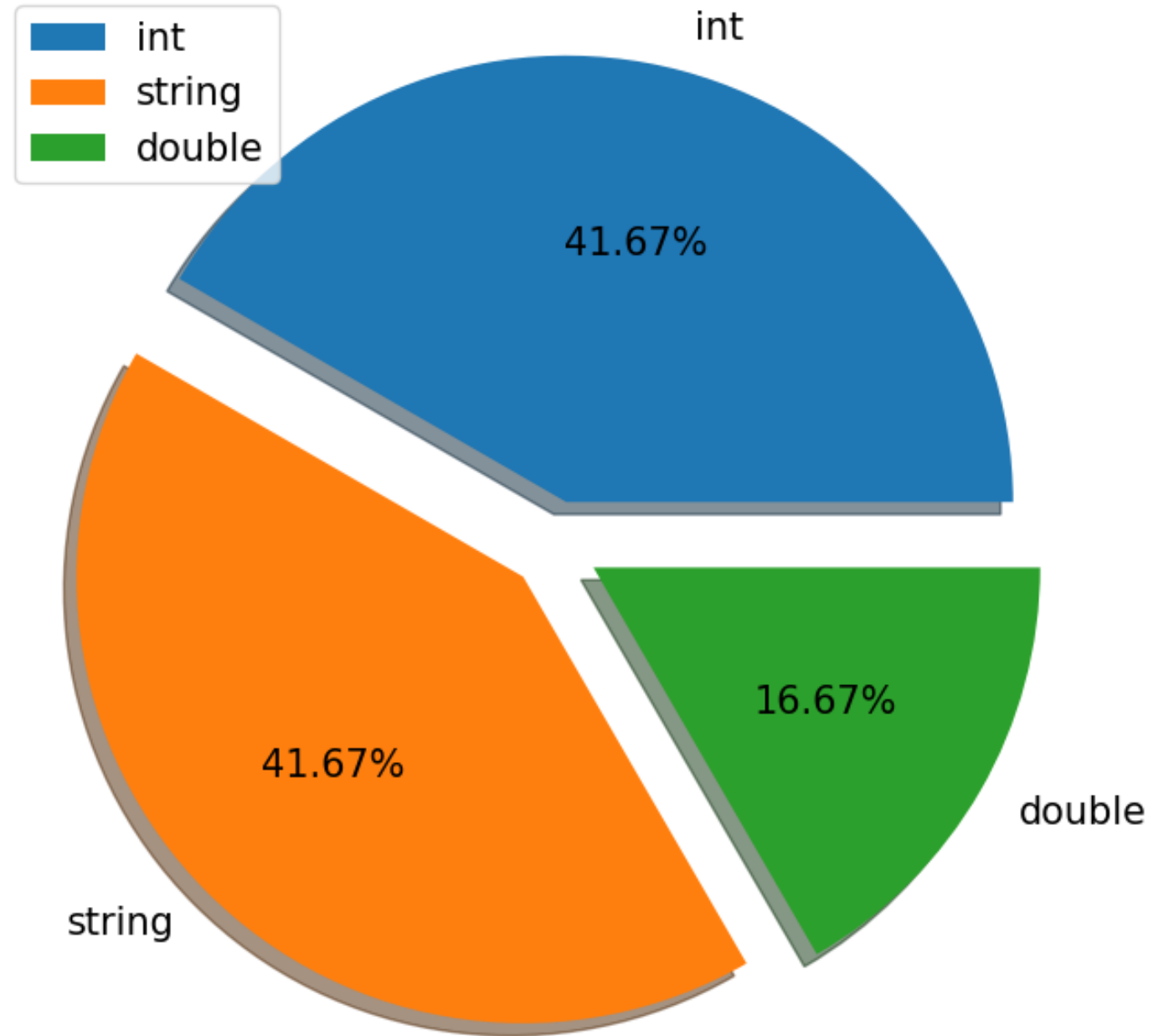




# DATA ANALYSIS & DATA VISUALIZATION

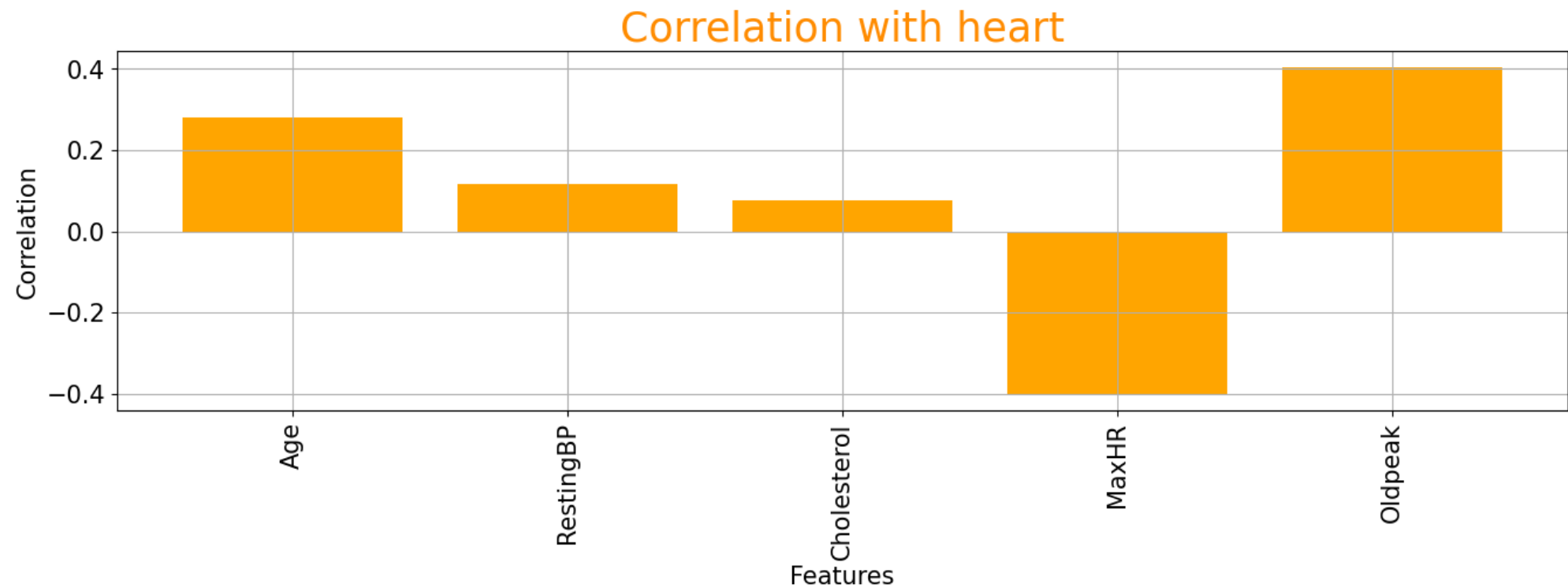


## Type of our data



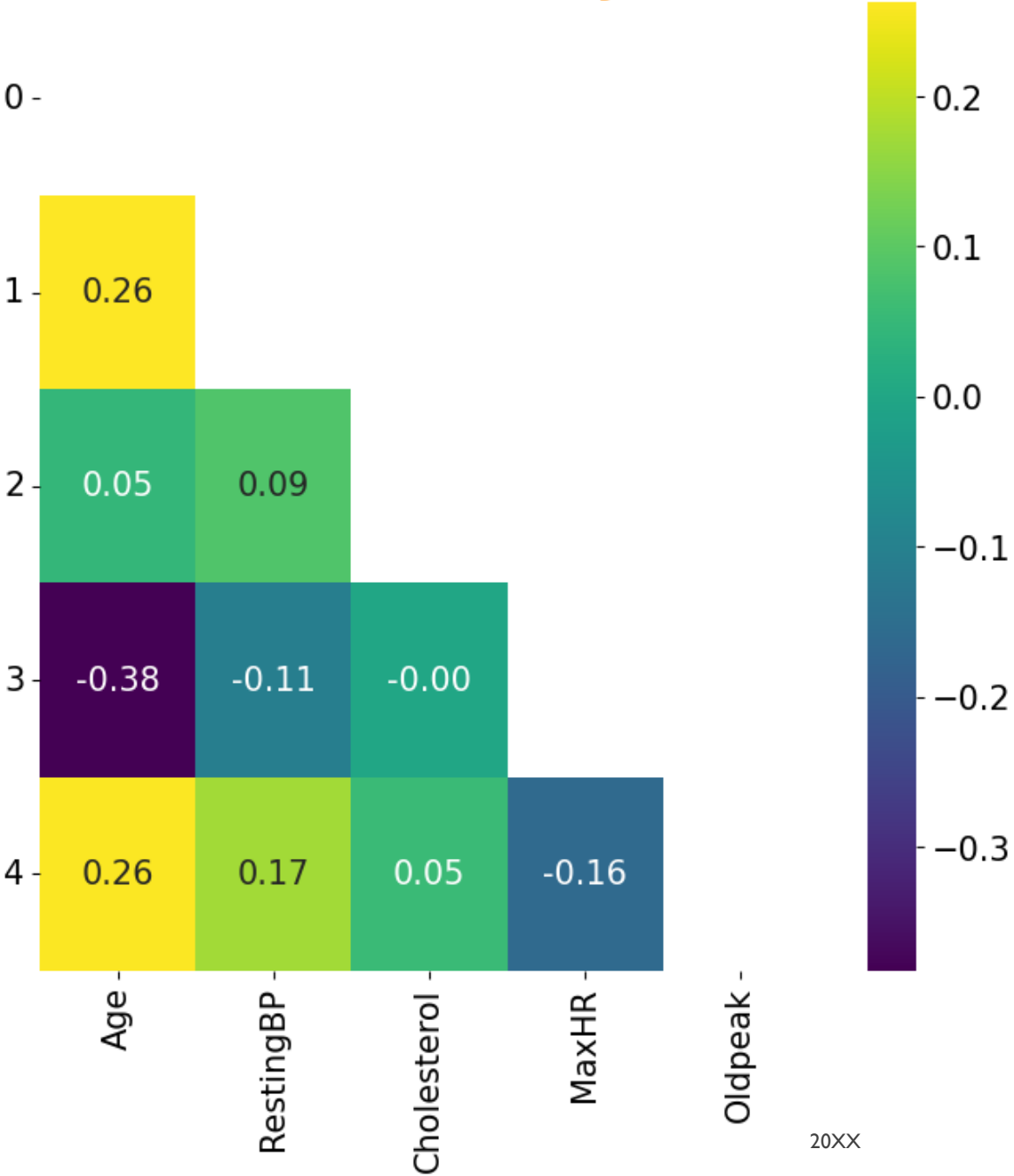
A PIE CHART FOR  
TYPE OF DATA

# A BAR PLOT FOR CORRELATION WITH HEART



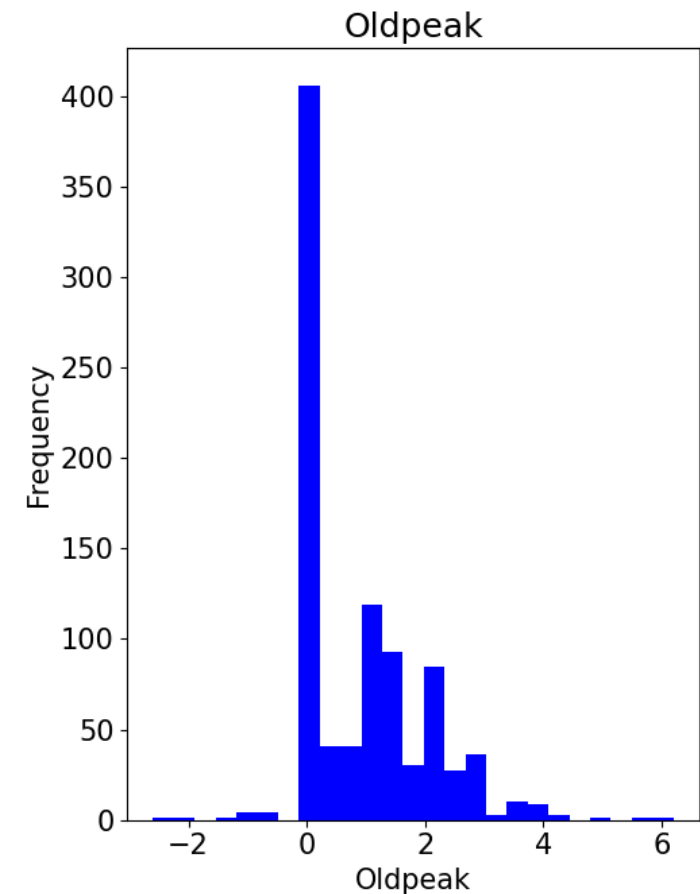
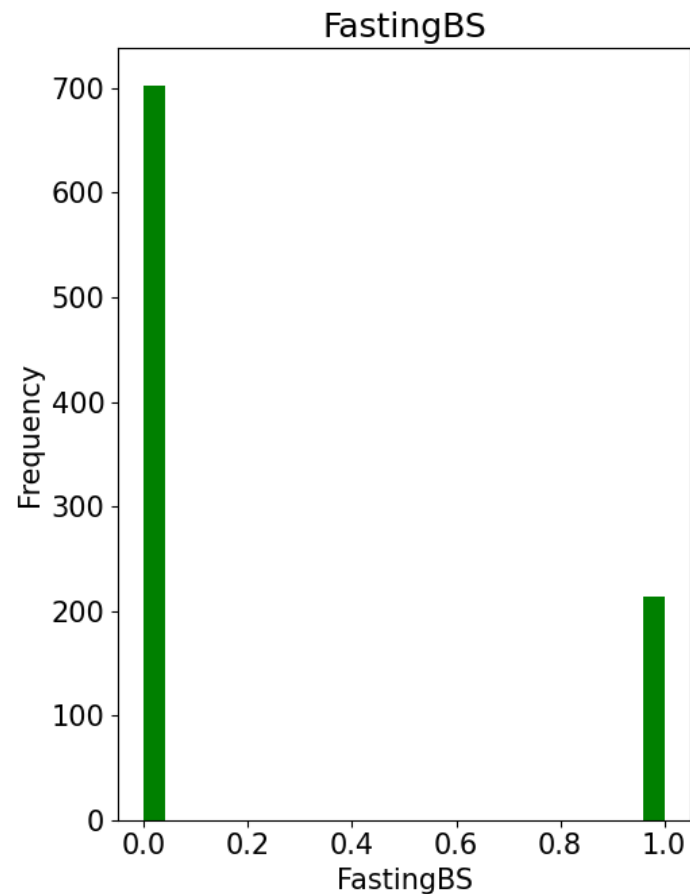
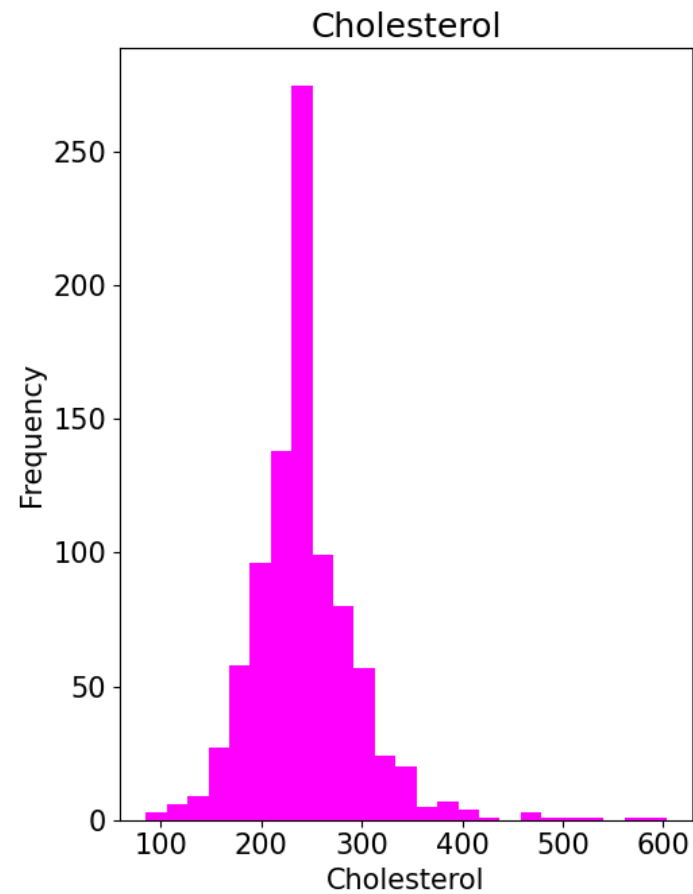
# Correlation Analysis

## HEATMAP FOR CORRELATION ANALYSIS

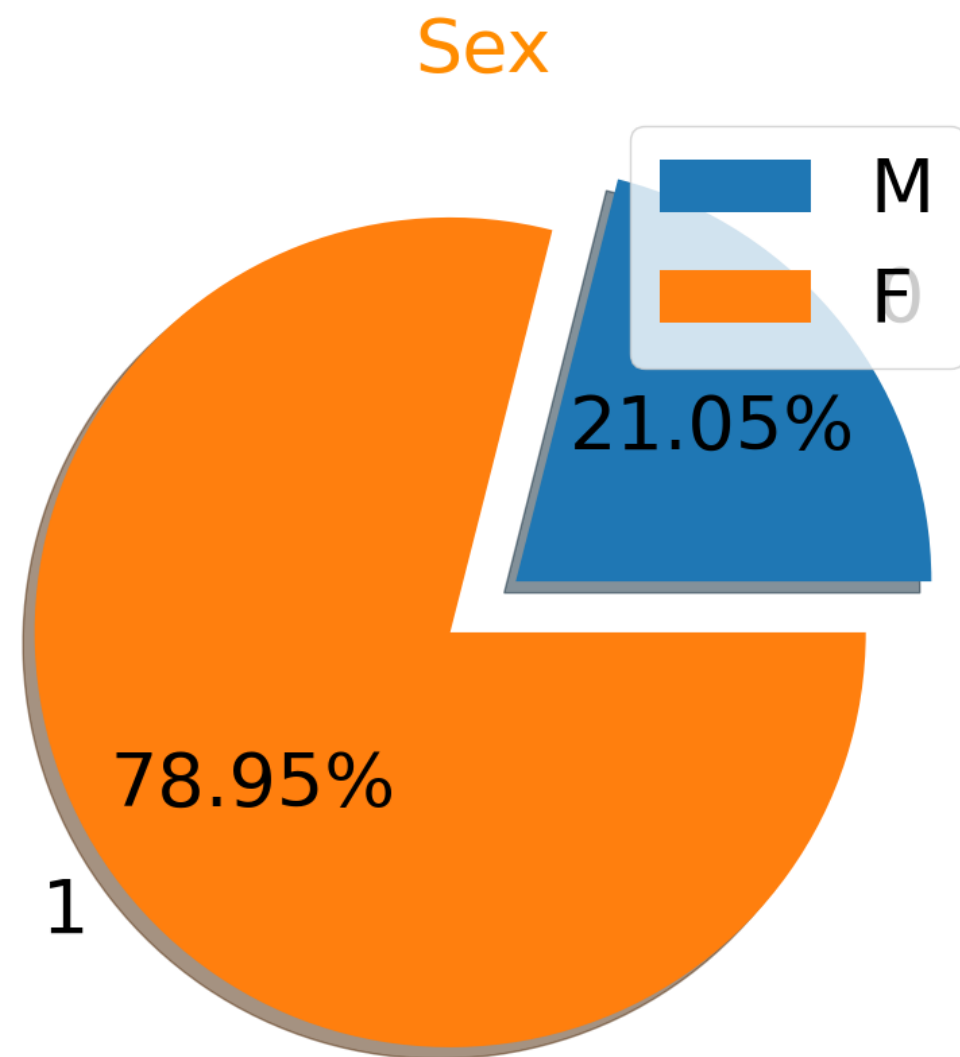
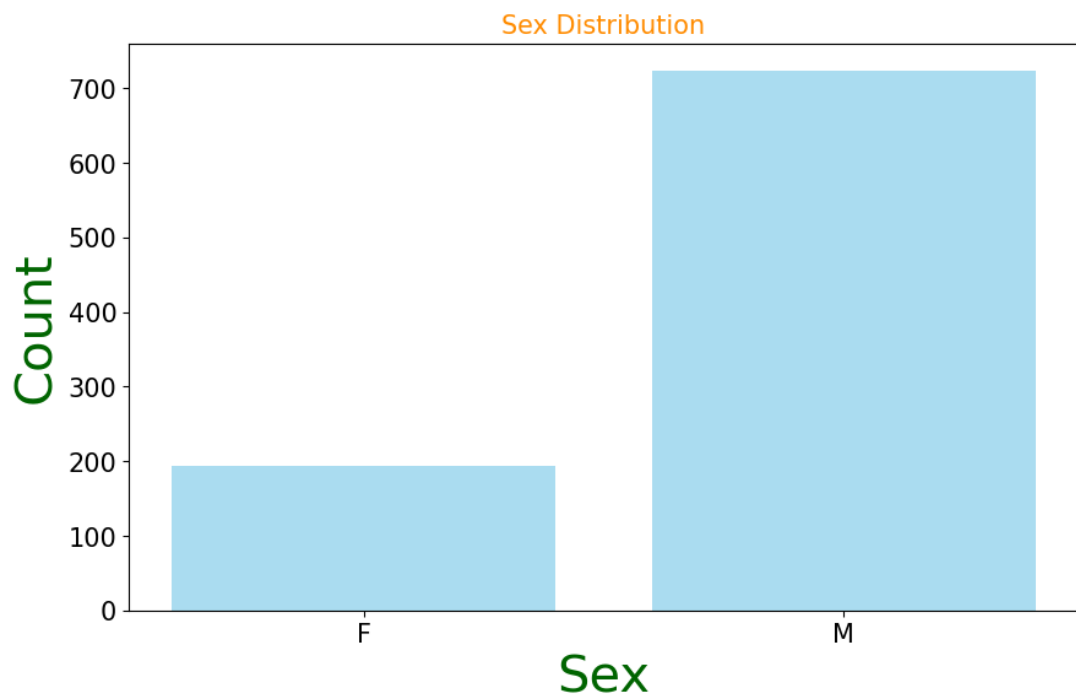


# HISTOGRAMS FOR THREE NUMERICAL COLUMNS

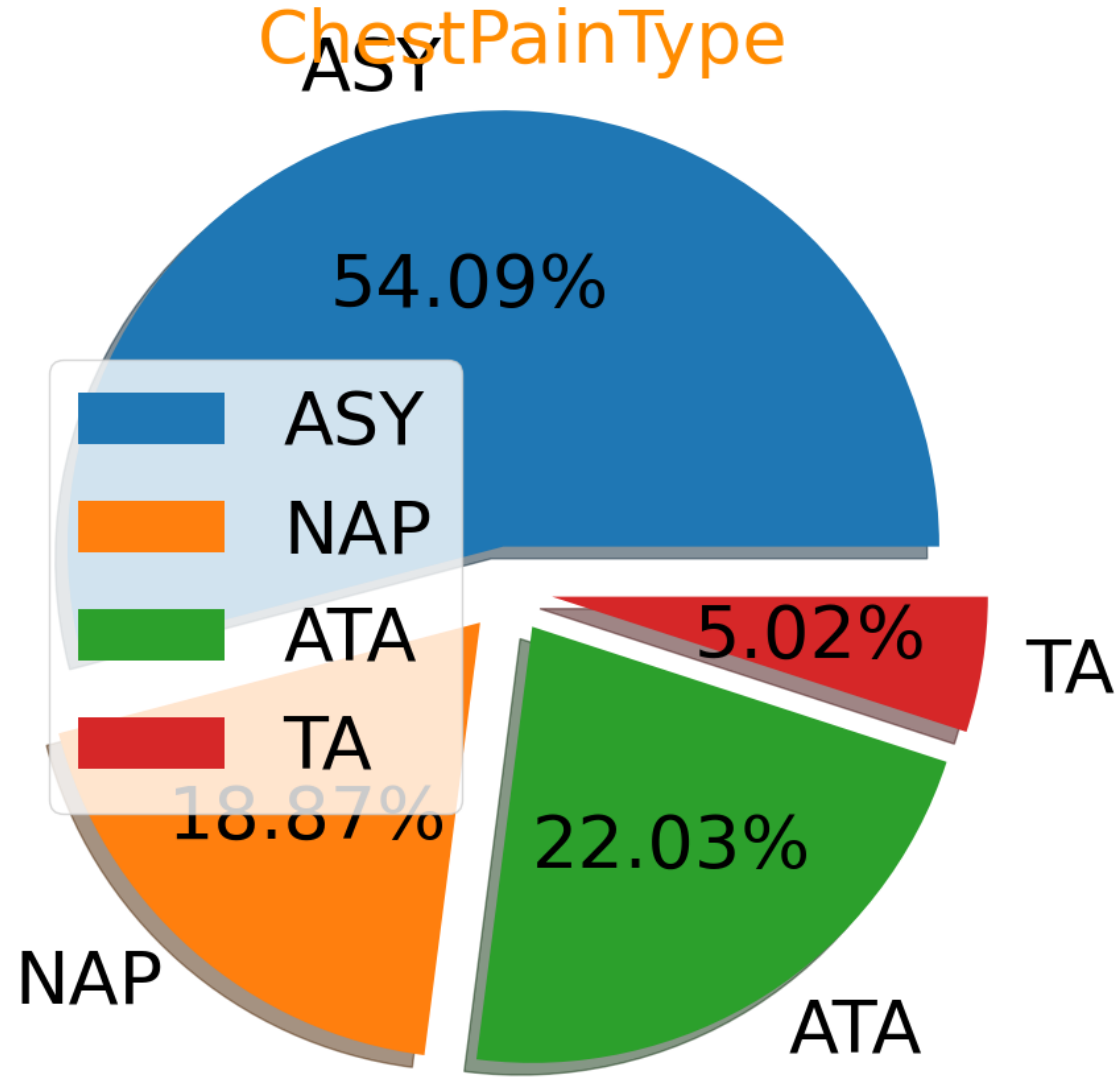
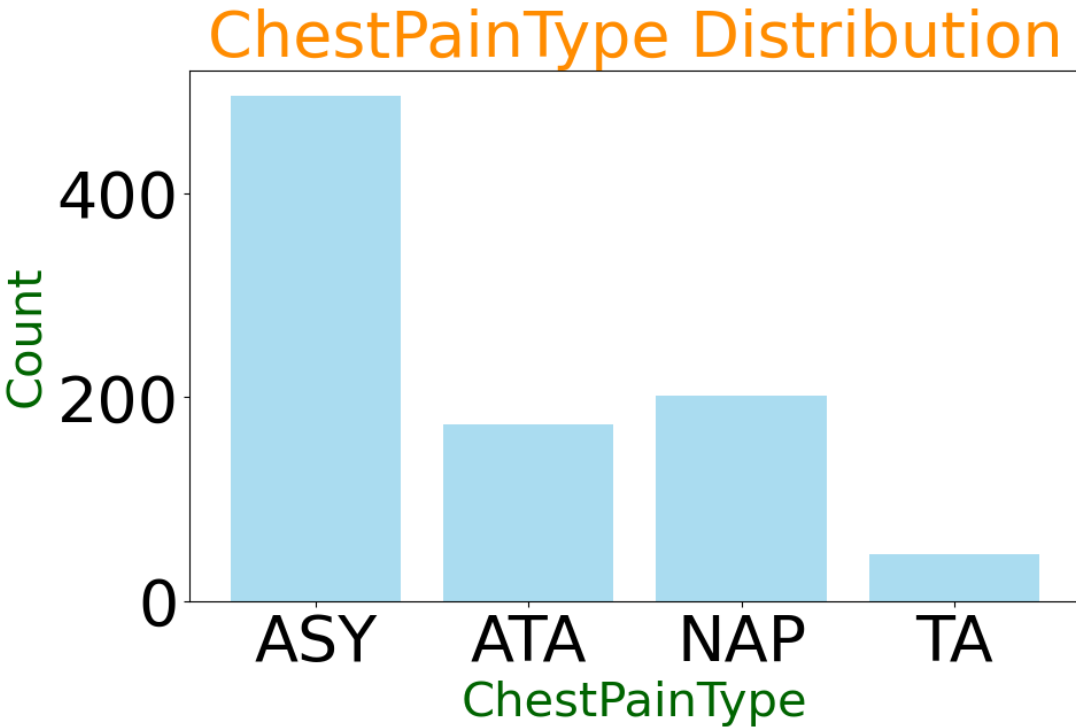
Subtitle



# VISUALIZATION OF SEX



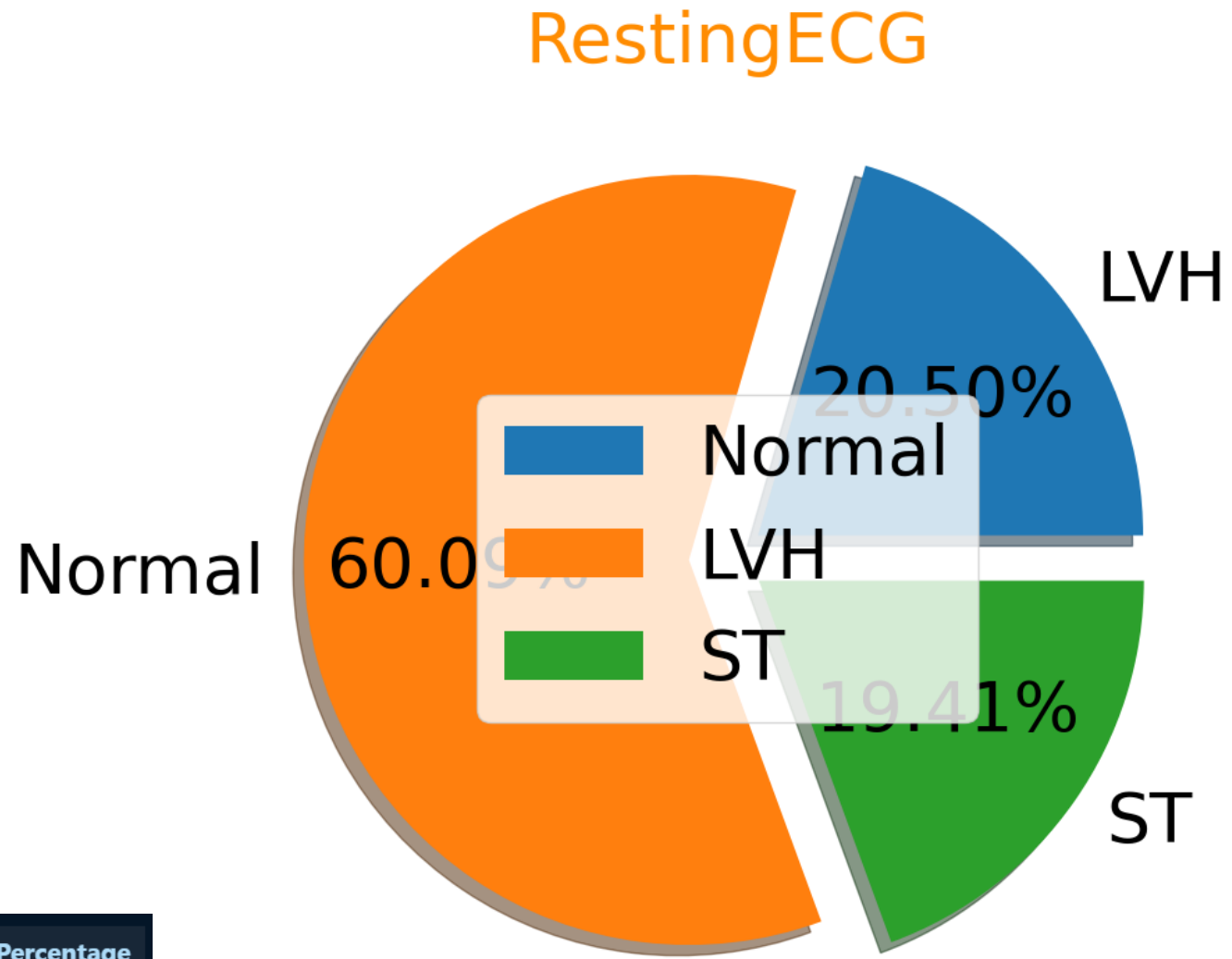
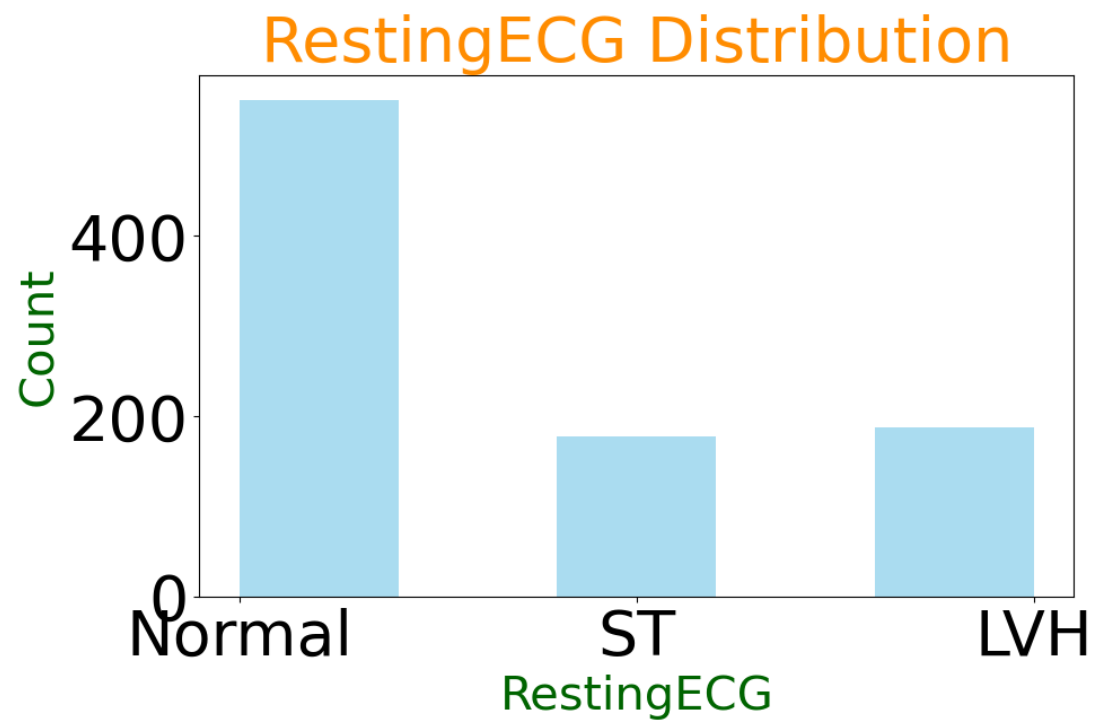
	Sex	Total	Percentage
0	F	193	21.05
1	M	724	78.95



VISUALIZATION OF  
CHEST PAIN TYPE

	ChestPainType	Total	Percentage
0	NAP	202	22.03
1	ATA	173	18.87
2	TA	46	5.02
3	ASY	496	54.09

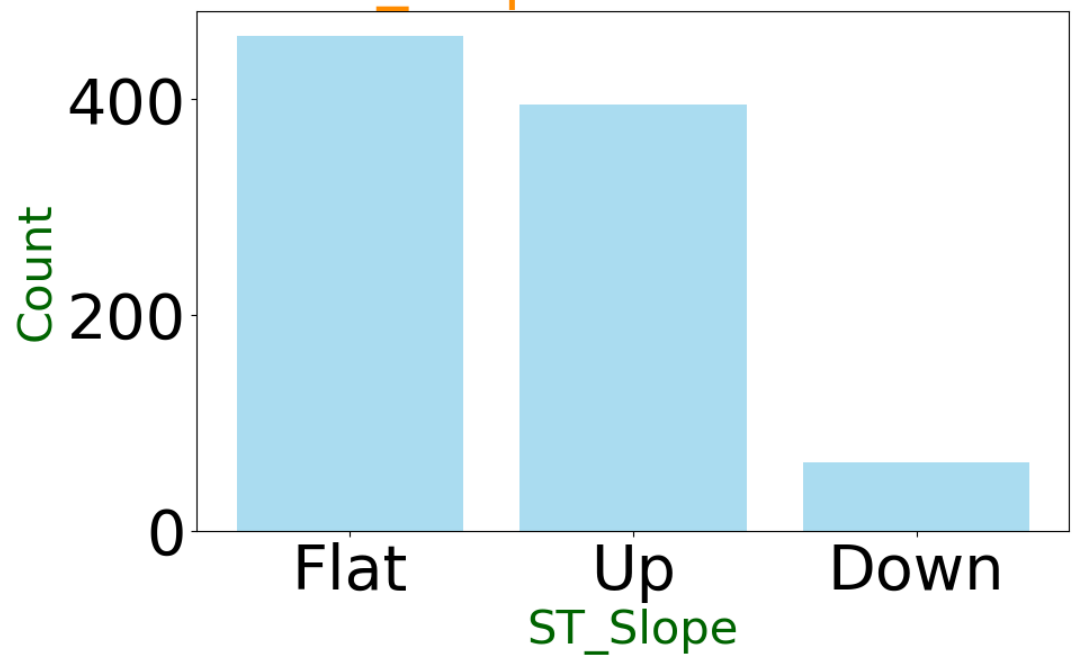




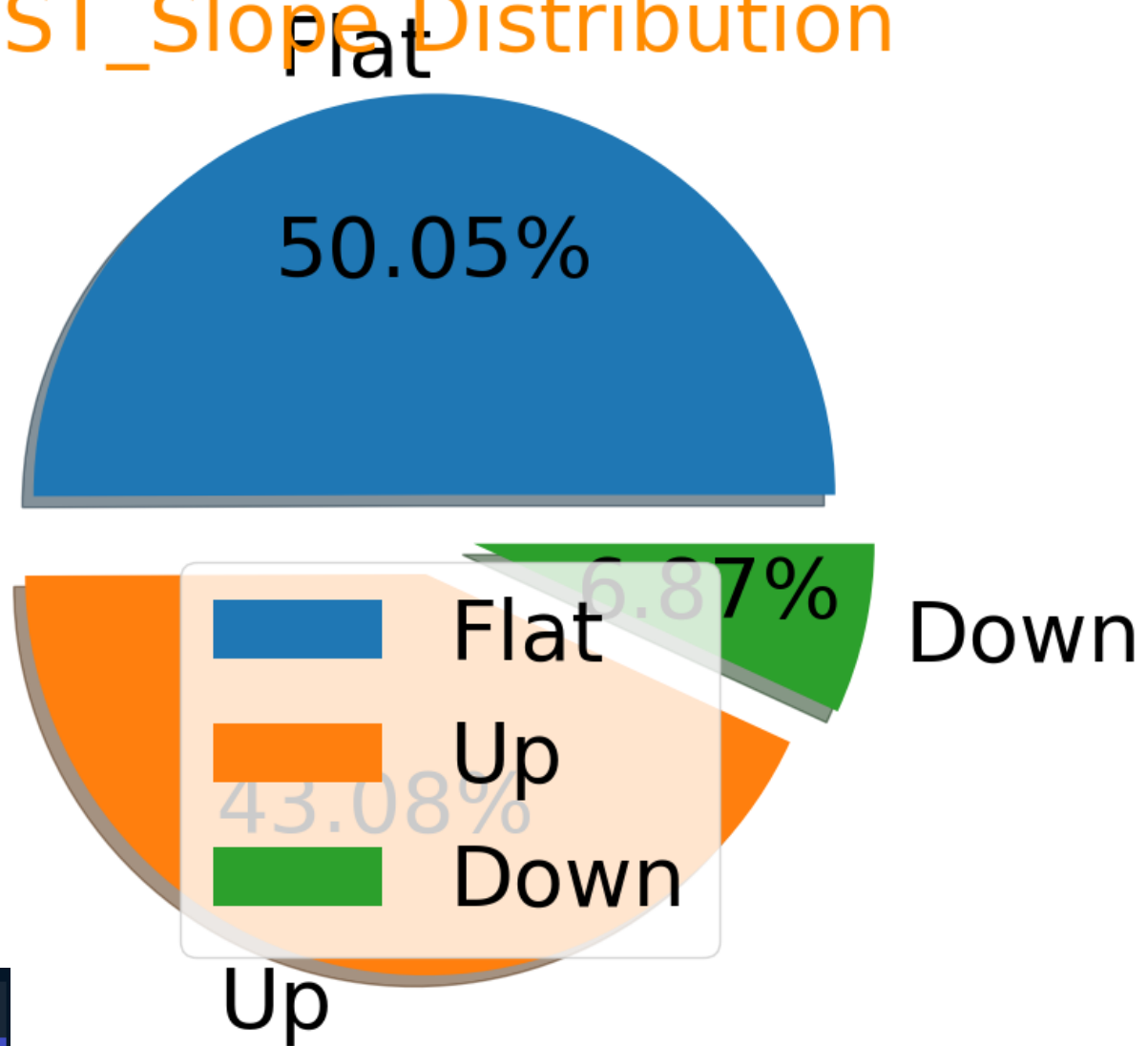
## VISUALIZATION OF RESTINGECG

	Total	Percentage
RestingECG		
Normal	551	60.09
LVH	188	20.50
ST	178	19.41

ST\_Slope Distribution

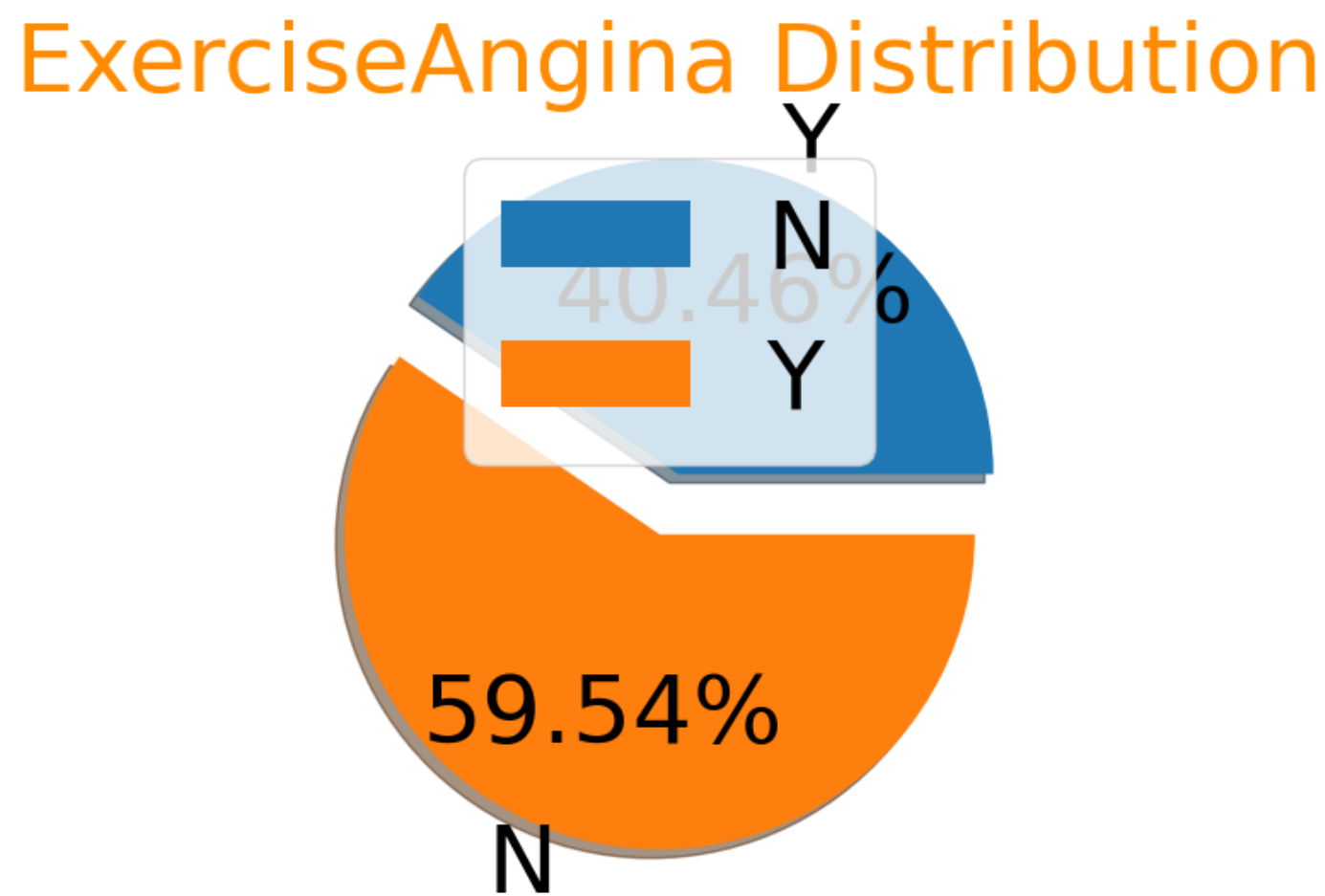
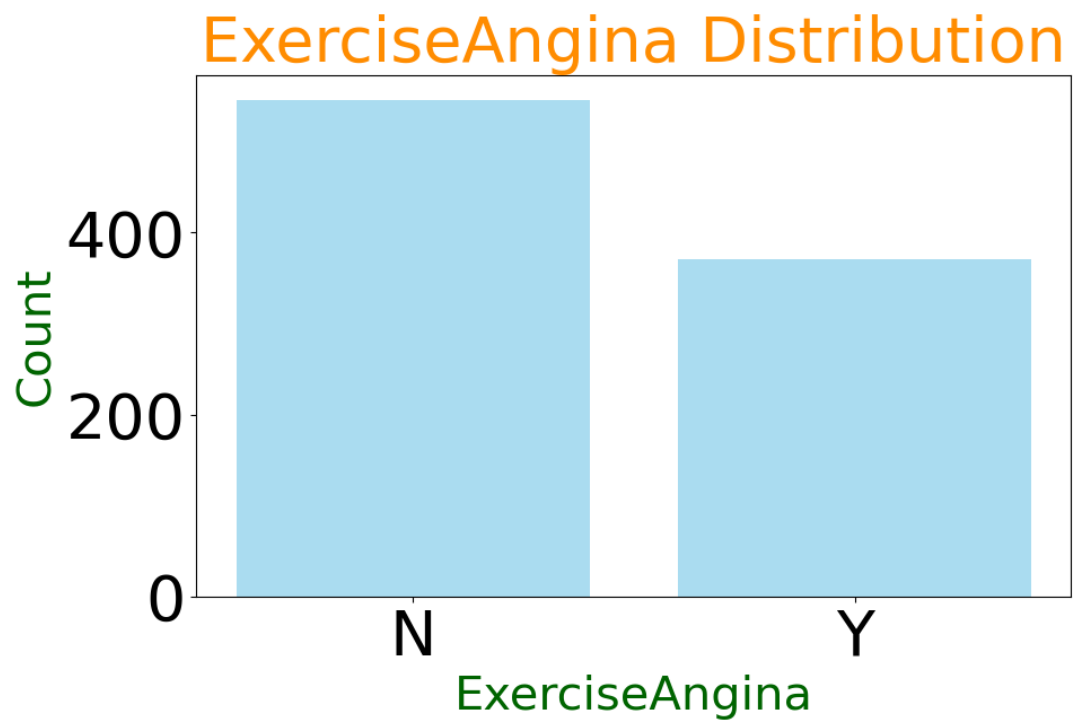


ST\_Slope Distribution



VISUALIZATION OF  
ST\_SLOPE

	Value	Total	Percentage
0	Y	371	40.460000
1	N	546	59.540000



VISUALIZATION OF  
EXERCISEANGINA

	Value	count	Percentage
0	Flat	459	50.050000
1	Up	395	43.080000
2	Down	63	6.870000

## RESULT OF ACCURACY SCORE

Model	Train Accuracy Score	Average Accuracy Score
Logistic Regression Model	84.45	91.35
Support Vector Machines Model	84.46	85.19
Random Forest Classifier Model	83.78	85.14
Multilayer Perceptron Classifier Model	81.08	79.48
Decision Tree Model	81.08	81.16



THANKS  
FOR  
LISTENING