

Table of contents

Introduction

About Cardiovascualr 01

Data Analysis

Exploratory data analysis about the data we have

Modeling

Reviewing the results of O3 O4 the models

Conclusions

Final conclusions in the insights and results we have



What is cardiovascular, and what are the diseases are related to it?

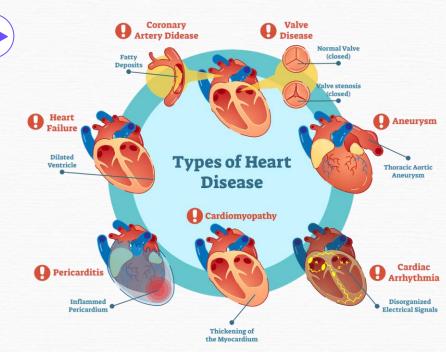


Cardiovascular:

 Cardiovascular disease (CVD) is a general term for conditions affecting the heart or blood vessels.

Cardiovascular Diseases:

- heart attacks
- Heart failure
- angina
- Strokes

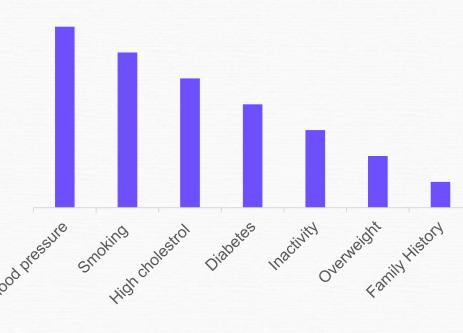


Most Common Causes for Cardiovascular diseases



Causes of cardiovascular diseases:

- Smoking
- Blood Pressure
- Family health history
- Inactivity
- Diabetes
- High cholesterol
- Being overweight



Here are the most common reasons are ranked based on NHS





Prevention:

- Stop smoking
- Health diet
- Exercise regularly
- Maintain healthy weight
- Medication

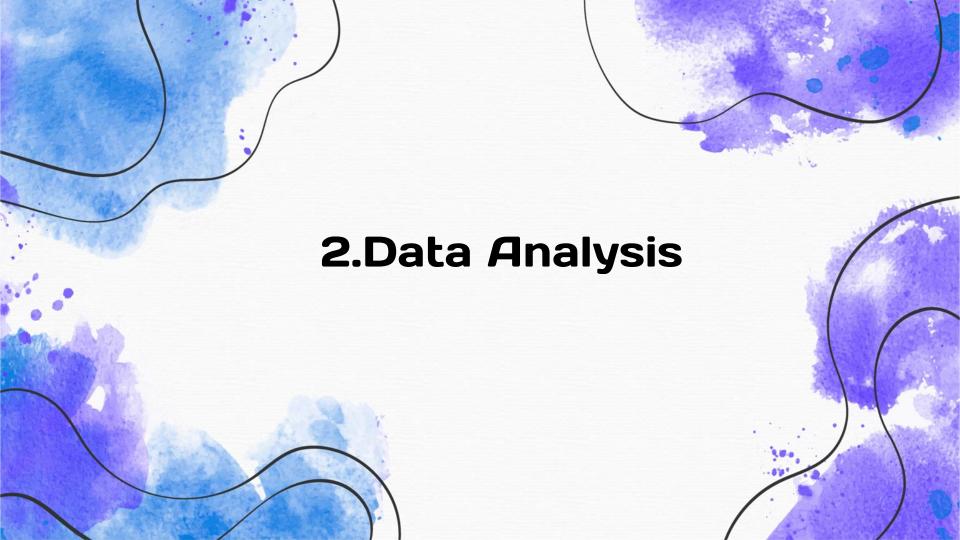
And so on based of each patient with his or her specific life style and what might be applicable treatment for them



There is a different way that might help preventing and help curing the disease before it might it be happening and that solution is:

Machine Learning Model

Machine learning will help using historical data of heart diseases to detect it's pattern and help us identify the diseases before it's actually happen, and start the treatment process early



Before diving in the machine learning process, we will do EDA

	age	gender	height	weight	ap_hi	ap_lo	cholest erol	gluc	smoke	alco	active	cardio	ВМІ
1	50.3	0	168	62.0	110	80	1	1	0	0	1	0	21.9
2	55.4	1	156	85.0	140	90	3	1	0	0	1	1	34.9
3	51.6	1	165	64.0	130	70	3	1	0	0	0	1	23.5
4	48.2	0	169	82.0	150	100	1	1	0	0	1	1	28.7
5	47.8	1	156	56.0	100	60	1	1	0	0	0	0	23.0

Features



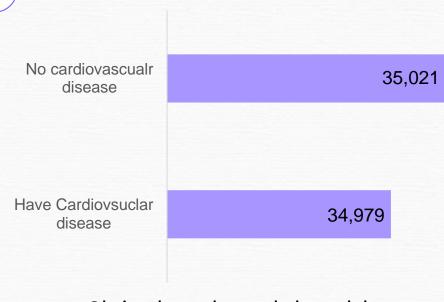
Target Variable

Patients without cardiovascular was higher than patients with it



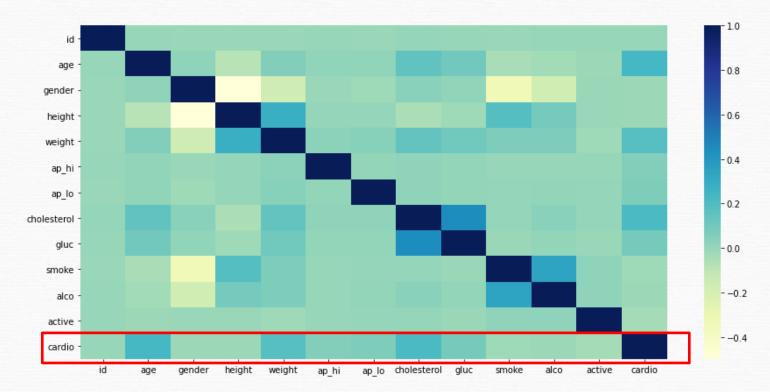
Cardiovascular disease (Our Target):

- Part of the analysis we do is to check number of patients and how many of them are facing diseases, so we can identify if the model and the data are balanced or not
- Also gain better knowledge about the patient data we have

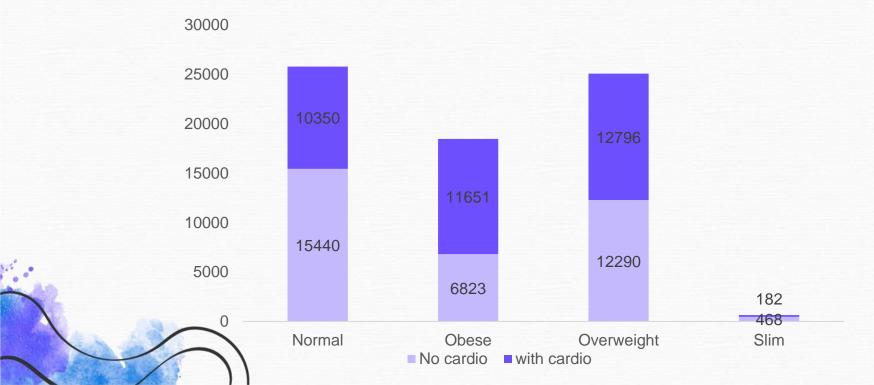


Obviously, we have a balanced dataset

4 Data elements that are positively correlated with cardiovascular



Below shows how patient who are obese and overweight have high cardiovascular disease





In the right section, we see a pie chart that is categorized for each male and female

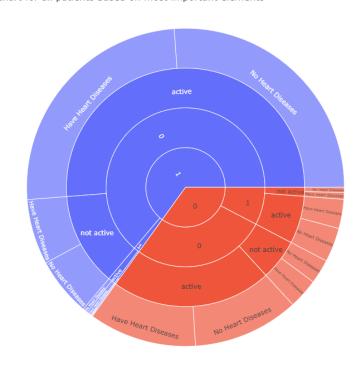
And it shows the characteristics of male and female patients and the relation of the data elements with cardiovascular

For example:

Female patients, who are not smoking and most of the time are active, most of them are not having heart diseases

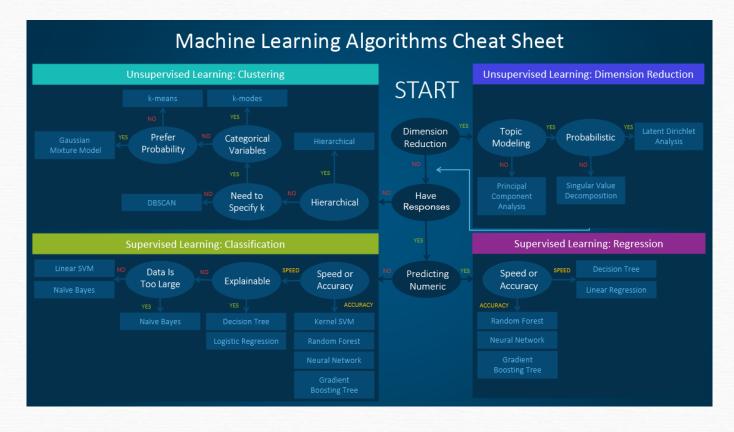
We can gain a lot of information from this chart

chart for all patients based on most important elements





Model Selection





Model Selection





- A multiple algorithms have been built as a based model to be compared and what model can be tuned, the following table shows the models

Model Name	Accuracy %			
Logistic Regression	72%			
KNN	66%			
Random Forrest	71%			
SVM	68%			
Decision Tree	63%			
Gradient Boosting	73%			

- The table shows that the logistic regression and Gradient boosting models, performed the best as a based models, so we will be choosing one of them and tune them

Focusing on the best accuracy





 Since we narrowed our models into two models, and they are logistic regression and gradient boosting

Model Name	Accuracy training	Accuracy validation
Logistic Regression	72%	72%
Gradient Boosting	73%	76%

When we compared both models, and checked the metrics we saw that gradient boosting algorithm performs better then logistic regression even with hyperparameter

Gradient Boosting Model





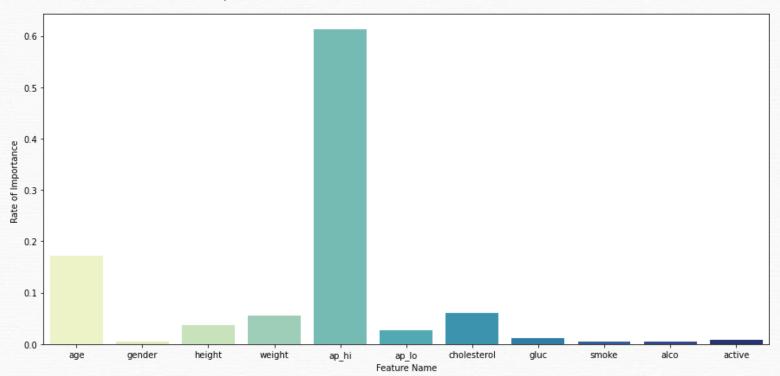
 We found our best fitted model for our case which is to identify patients with cardiovascular based on patient information

Model Name	Accuracy training	Accuracy validation	Accuracy Testing
Gradient Boosting	73%	76%	74%

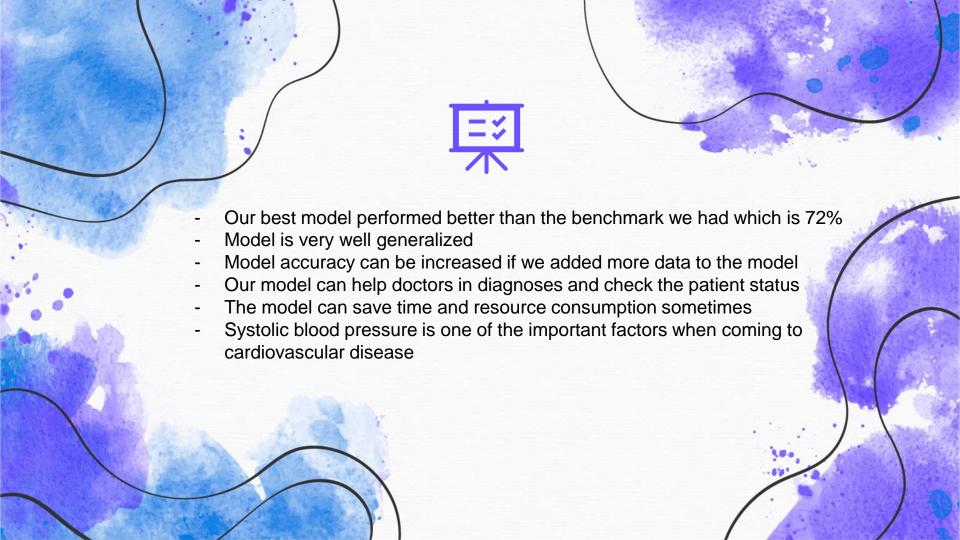
Lastly we predicted the results using our testing data and from the accuracy in general we can say that the model is very well generalized and learned good and not overfitting

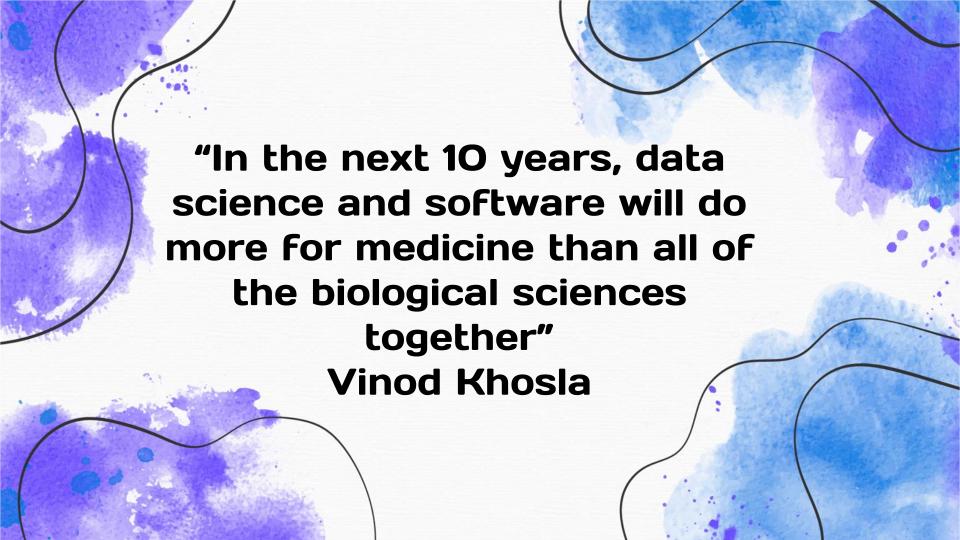
Feature Importance

From below we see that Systolic blood pressure (ap_hi) is the most important feature out of the total elements we have









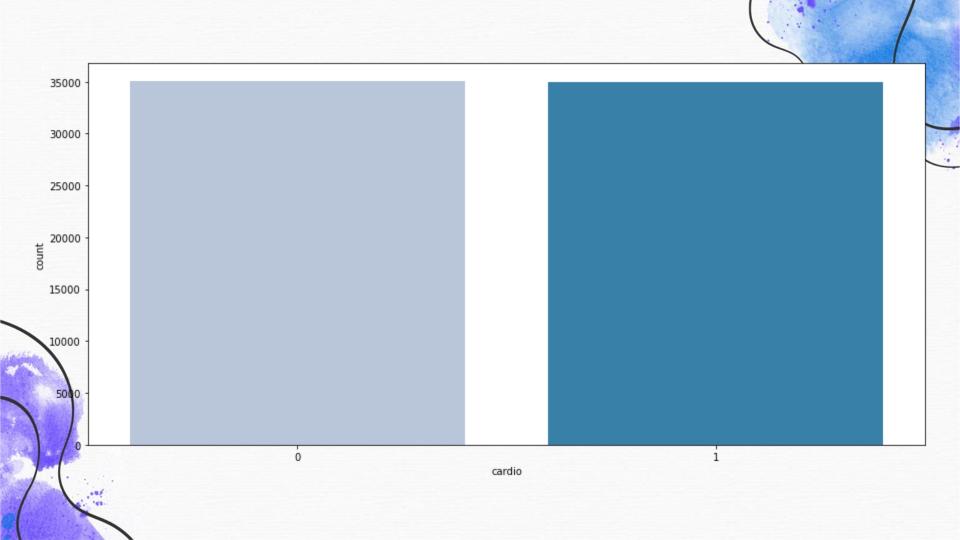


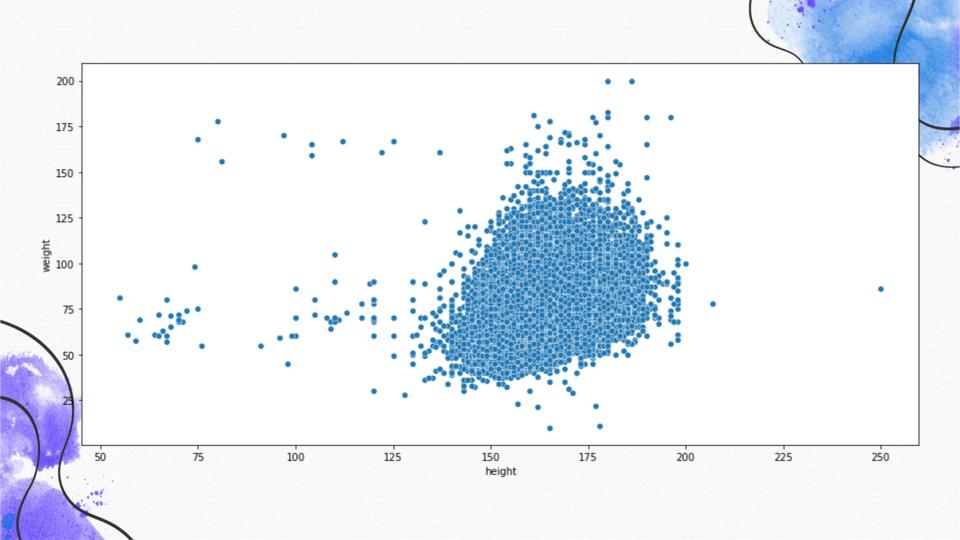
Research resources

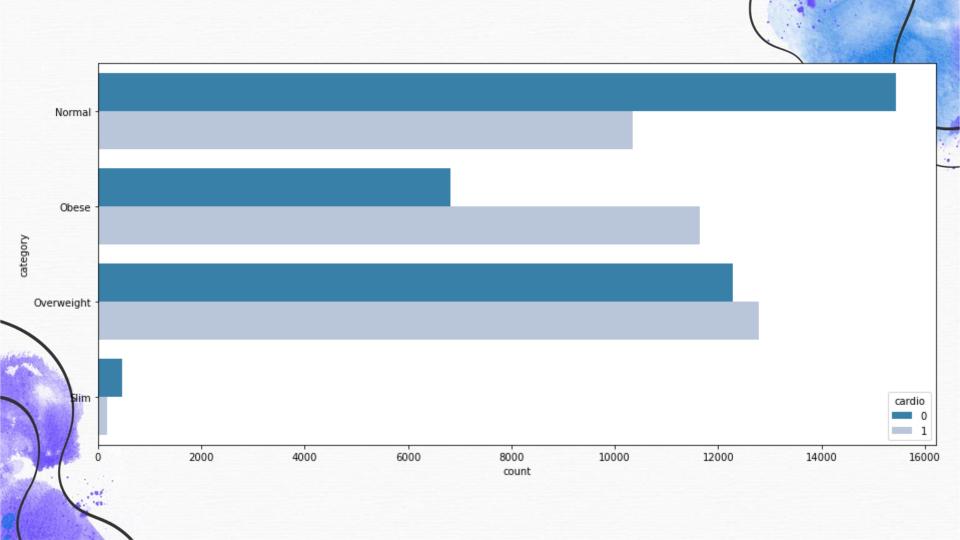
- https://quotefancy.com/quote/1272586/Vinod-Khosla-In-the-next-10-years-data-science-and-software-will-domore-for-medicine
- https://www.kaggle.com/sulianova/cardiovasculardisease-dataset
- https://www.nhs.uk/conditions/cardiovascular-disease/
- https://plotly.com/python/sunburst-charts/
- https://machinelearningmastery.com/calculate-featureimportance-with-python/





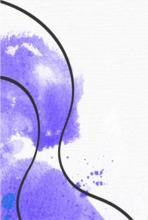


















#348ce7