

- ① Imported csv file using pandas read-csv
- ② checked shape of both DataFrames. using shape function
- ③ checked Null values using isnull function
- ④ Dropped some ^{duplicated} rows and columns which aren't required. used drop function using value-count
- ⑤ checked unique values from categorical columns, & removed outliers values with deviate signifi. from rest.
 ⇒ used boxplot to check outliers in Numerical columns & removed outliers (if exists)
- ⑥ Visualised data using different graphs (Pie, Bar graph, Histogram, Scatter Plot) using Matplotlib & Seaborn library.
- ⑦ Used Label Encoder to convert categorical columns to numerical columns.
- ⑧ Selected Features & Targeted columns to numerical columns.
- ⑨ Created Training & testing Dataset using train-test split.
- ⑩ Used ~~for~~ classification model & checked accuracy, confusion matrix.
 1. LR
 2. Decision Tree classification
 3. Random Forest
 4. SVM classification (Support Vector Machine)
 5. KNN (K Near Neighbour)
 6. XGBoost classification.
- 75% → Accuracy, precision score isn't good.
- ⑪ Normalized data using MinMaxScaler,
- ⑫ Again the data was imbalanced, used ~~imbalan~~ imblear over-sampling to balance data

Again used all ML model, and got max. accuracy
84-14% using XGBoost Model.

② used k-fold & stratified shuffle split for validation.

pickle →

Std. way of serializing objects in python.

↓

converting objects into stream of bytes. to

store the object or transmit it to memory or a file.

Reg.) class

↓

Regression algo → used to pred

Confusion Matrix

↳ It's $N \times N$ matrix used to evaluate the performance of classification model, when N is number of target values.

		Actual values	
		+ve	-ve
Predicted values	+ve	TP	FP
	-ve	FN	TN

Sick ppl are correctly predicted by sick model

Healthy people are incorrectly predicted by sick model

Sick ppl are incorrectly predicted as not sick by the model

Healthy ppl correctly predicted as not sick by the model

12 features

↓

30 features