

THYROID DISEASE DETECTION

Detailed Project Report

Kunal Aggarwal
Prachi Bindal

INTRODUCTION

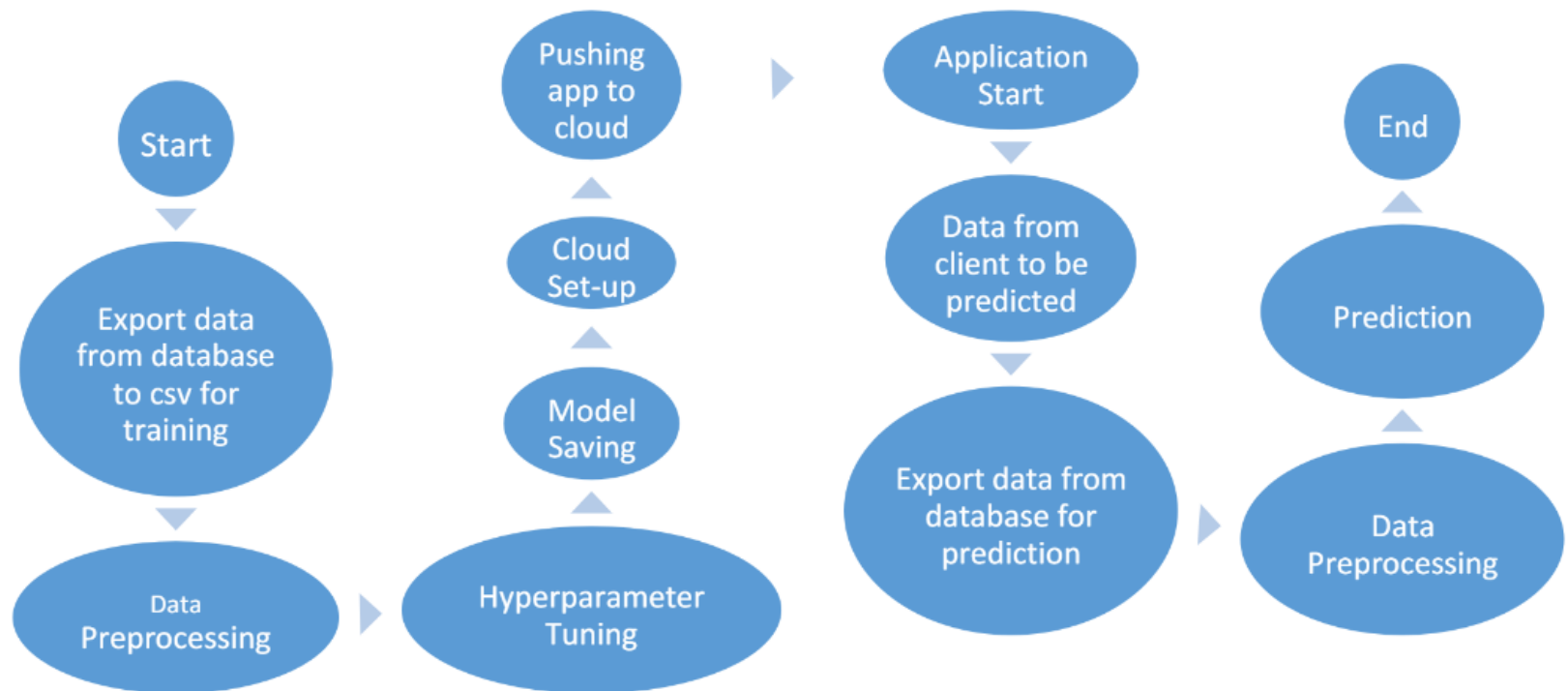
At least a person out of ten is suffered from thyroid disease in India. The disorder of thyroid disease primarily happens in the women having the age of 17–54. The extreme stage of thyroid results in cardiovascular complications, increase in blood pressure, maximizes the cholesterol level, depression and decreased fertility. The hormones, **total serum thyroxin (T4)** and **total serum triiodothyronine (T3)** are the two active thyroid hormones produced by the thyroid gland to control the metabolism of body. For the functioning of each cell and each tissue and organ in a right way, in overall energy yield and regulation and to generate proteins in the ordnance of body temperature, these hormones are necessary.

Hyperthyroidism and **Hypothyroidism** are the most two common diseases caused by irregular function of thyroid gland. Thyroid disorder can speed up or slow down the metabolism of the body. In the world of rising new technology and innovation, health care industry is advancing with the role of Artificial Intelligence. Machine learning algorithms can help to early detection of the disease and to improve the quality of the life. This study demonstrates the how different classification algorithms can forecasts the presence of the disease. Different classification algorithms such as Logistic regression, Random Forest, Decision Tree, Naïve Bayes, Support Vector Machine, XG Boost, KNN have been tested and compared to predict the better outcome of the model.

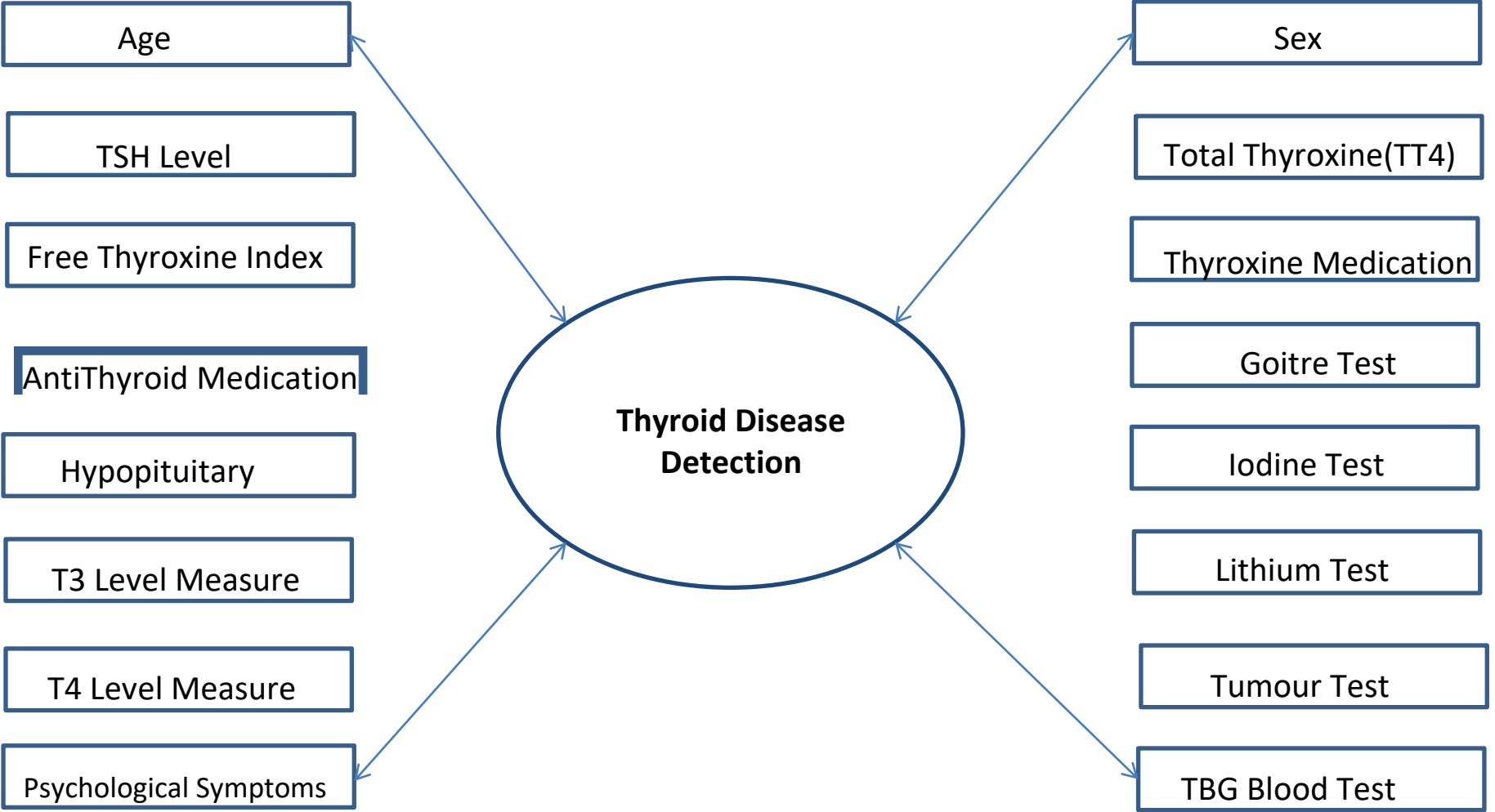
OBJECTIVE

The main goal of this project is to predict the risk of hyperthyroid and hypothyroid based on various factors of individuals. Thyroid disease is a common cause of medical diagnosis and prediction, with an on set that is difficult to forecast in medical research. It will play a decisive role in order to early detection, accurate identification of the disease and helps the doctors to make proper decisions and better treatment.

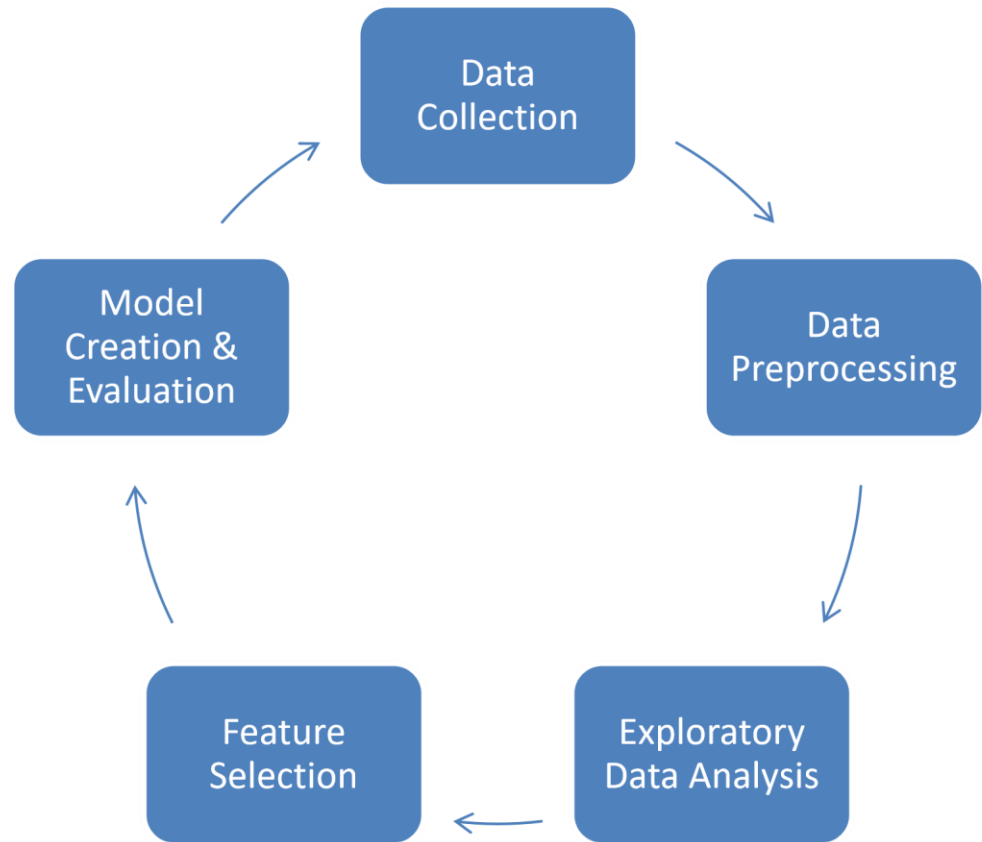
ARCHITECTURE



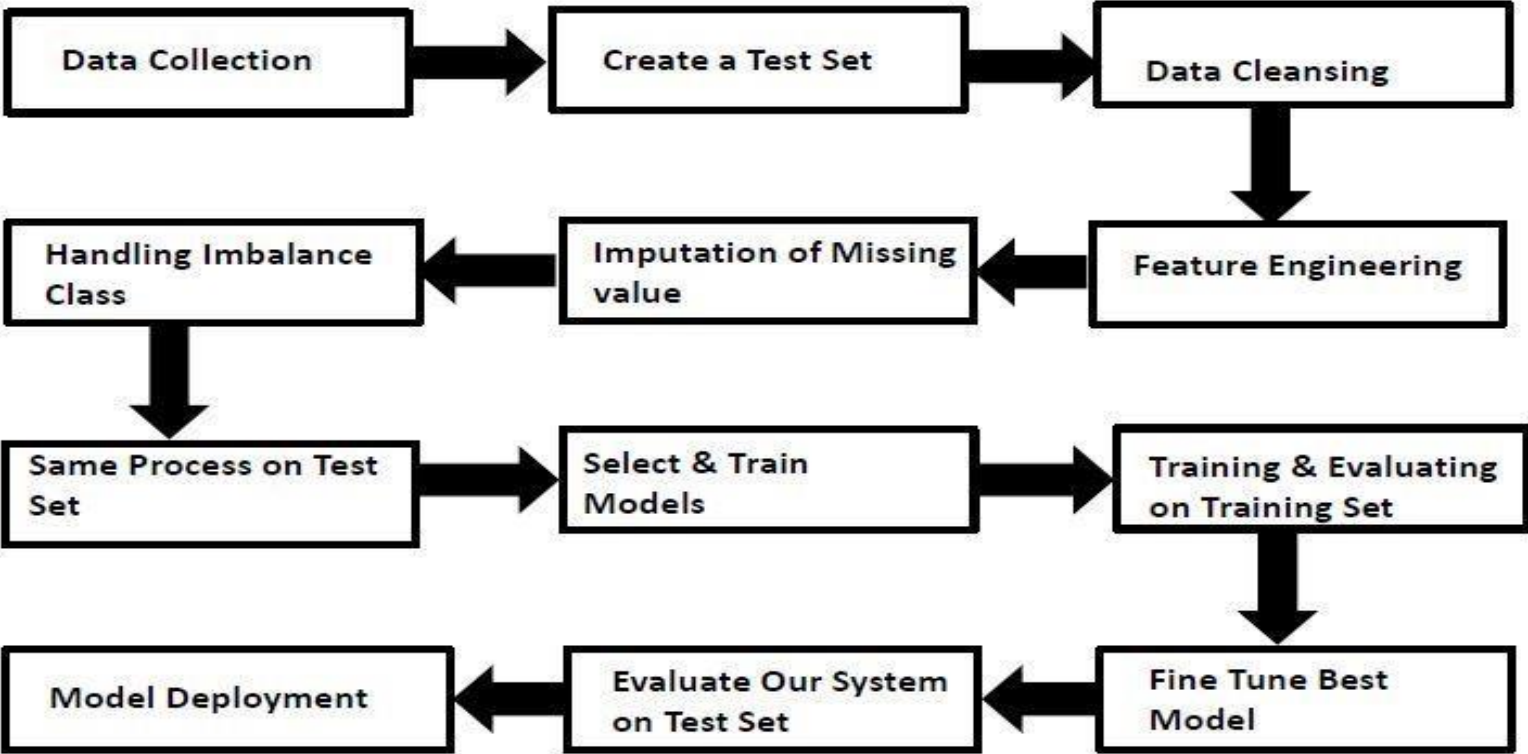
DATASET



Data Analysis Steps



MODEL TRAINING AND VALIDATION WORKFLOW



MODEL TRAINING AND VALIDATION WORKFLOW

Data Collection

- Thyroid Disease Data Set from UCI Machine Learning Repository
- For Data Set: <https://archive.ics.uci.edu/ml/datasets/thyroid+disease>

Data Pre-Processing

- Missing values handling by Simple imputation (Used KNN Imputer)
- Outliers' detection and removal by boxplot and percentile methods
- Categorical features handling by ordinal encoding and label encoding
- Feature scaling done by Standard Scalar method
- Imbalanced dataset handled by SMOTE -Over sampling
- Drop unnecessary columns

MODEL TRAINING AND VALIDATION WORKFLOW

Model Creation and Evaluation

- Various classification algorithms like Random Forest, XG Boost, KNN, etc. tested.
- Random Forest, XGBoost and KNN all were given better results. Random Forest was chosen for the final model training and testing.
- Hyper parameter tuning was performed.
- Model performance evaluated based on accuracy, confusion matrix, classification report.

Random Forest Classifier Model

INTRODUCTION

It is a decision-tree-based ensemble Machine Learning algorithm which combines the output of multiple decision trees to reach a single result.

The Random Forest Classifier is a supervised learning algorithm which we can use for regression and classification problems. It is among the most popular machine learning algorithms comes under bagging ensemble technique.

Random Forest Classifier being ensemble algorithm tends to give more accurate result. This is because it works on the principle i.e., it creates the random forest by combining N decision tree, and make predictions for each tree created in the first phase. Even if one or few decision tree are prone to noise, overall results would tend to be correct.

Reason to use Random Forest Classifier model:

- It takes less training time as compared to other algorithms.
- It gives better model performance.

DATABASE CONNECTION & DEPLOYMENT

Database Connection

- MongoDB Database is used for this project.

thyroidDetection.patients

81DOCUMENTSINDEXES

DocumentsAggregationsSchemaExplain PlanIndexesValidation

FilterType a query: { field: 'value' }

ResetFindMore Options

ADD DATAEXPORT COLLECTION

1 - 8 of 8

_id: ObjectId('648dc830cd228def9136e91e')

age: 24

sex: 1

TSH: 456

TT4: 34

FTI: 45

T3: '4'

T4U: 1

on_thyroxine: 1

on_antithyroid_medication: 1

goitre: 0

hypopituitary: 0

psych: 1

lithium: 0

TSH_measured: 1

TT4_measured: 1

T4U_measured: 1

T3_measured: 1

query_on_thyroxine: 1

query_hypothyroid: 1

query_hyperthyroid: 1

T131: 1

thyroid_surgery: 1

pregnant: 0

sick: 1

SHOW 2 MORE FIELDS

Model Deployment

- The final model is deployed on Render using Flask framework.

The word "render" is displayed in a large, bold, dark blue sans-serif font. A thin horizontal line is positioned above the text.

FREQUENTLY ASKED QUESTIONS

Q1) What is the source of data?

The data for training is obtained from famous machine learning repository.

UCI Machine Learning Repository: <https://archive.ics.uci.edu/ml/datasets/thyroid+disease>

Q2) What was the type of data?

The data was the combination of numerical and Categorical values.

Q3) What's the complete flow you followed in this Project?

Refer slide 7th, 8th and 9th for better understanding.

Q4) After the File validation what you do with incompatible file or files which didn't pass the validation?

Files like these are moved to the Archive Folder and a list of these files has been shared with the client and we removed the bad data folder.

Q 6) What techniques were you using for data pre-processing?

- Removing unwanted attributes

- Visualizing relation of independent variables with each other and output variables
- Checking and changing Distribution of continuous values
- Removing outliers
- Cleaning data and imputing if null values are present.
- Converting categorical data into numeric values.

Q 7) How training was done or what models were used?

- First Data validation done on raw data and then good data insertion happen in DB.
- Then Data preprocessing done on final CSV file received from DB.
- Various model such as Decision Tree, Random Forest and XGBoost models are trained and based on performance, model is saved.

Q 8) How Prediction was done?

- The client fills the required inputs which is visible on the homepage of API.
- After filling all the required inputs, prediction was made and client sees the desired result.

Q 9) What are the different stages of deployment?

- After model training and finalizing all models. We created required files for deployment.
- Finally deployed our model over a cloud platform named as **Render**.

Q 10) How is the User Interface present for this project?

- For this project we have made only one type of UI.
- It is for one user input prediction.
- It very user friendly and easy to use.

THANK YOU