

THYROID DISEASE DETECTION

Detailed Project Report

Manish Kumawat
Data Science Intern at
PwSkills

INTRODUCTION

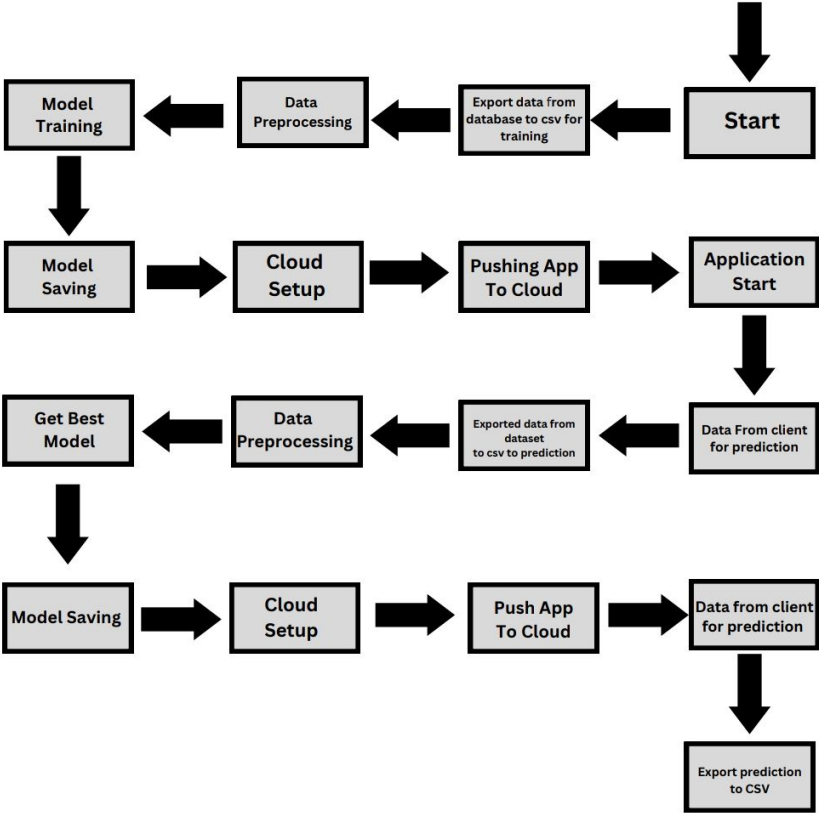
At least a person out of ten is suffered from thyroid disease in India. The disorder of thyroid disease primarily happens in the women having the age of 17–54. The extreme stage of thyroid results in cardiovascular complications, increase in blood pressure, maximizes the cholesterol level, depression 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, Support Vector Machine, KNN have been tested and compared to predict the better outcome of the model.

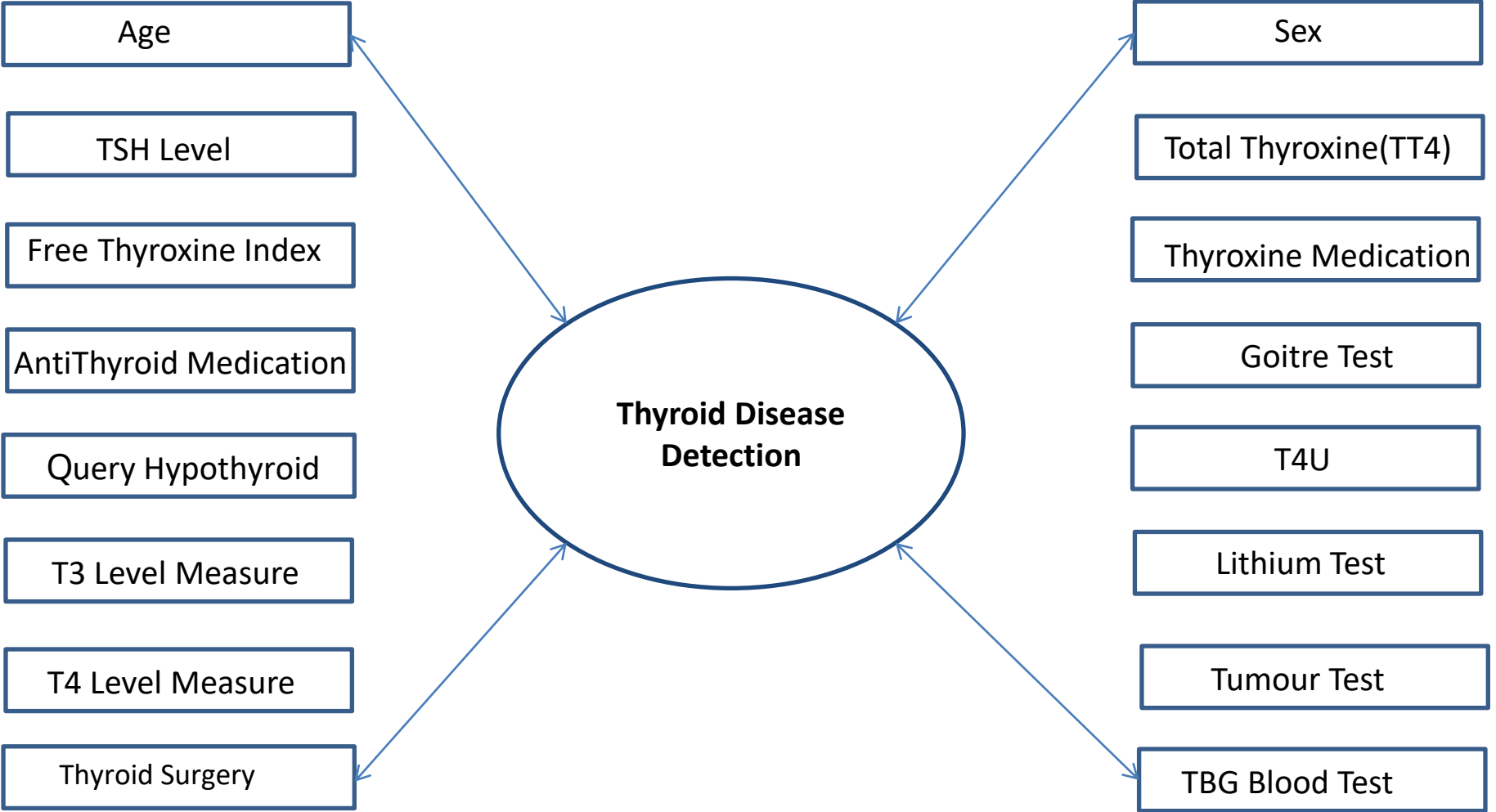
OBJECTIVE

The main goal of this project is to predict the risk of 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 fore cast 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.

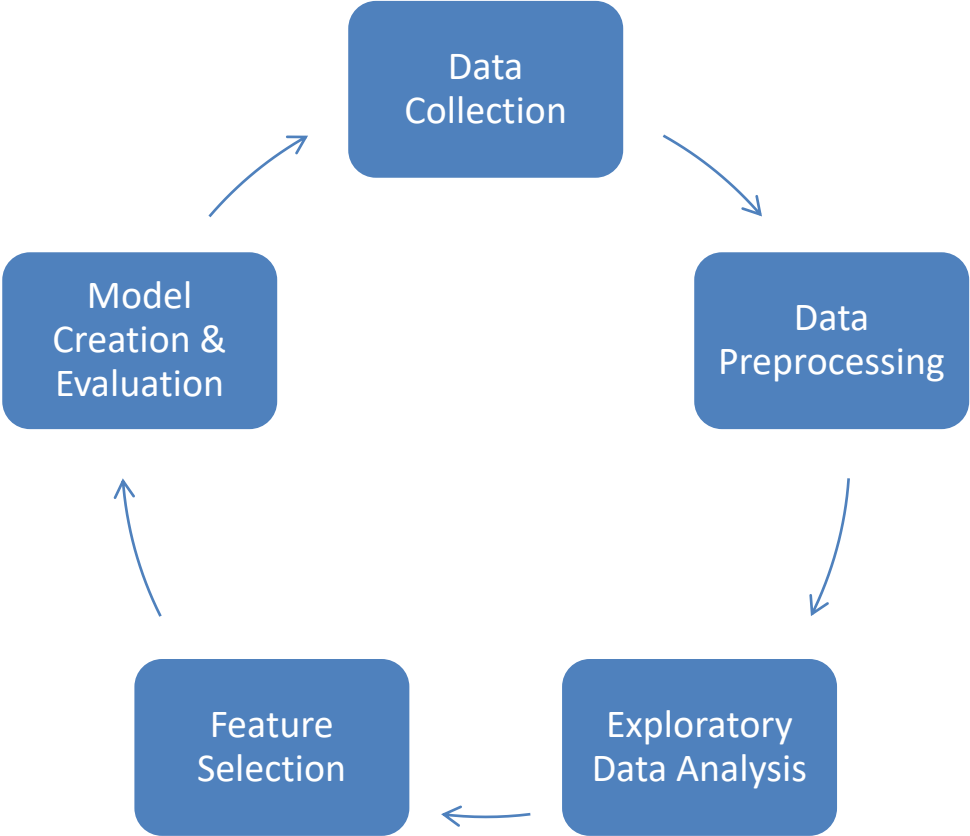
ARCHITECTURE



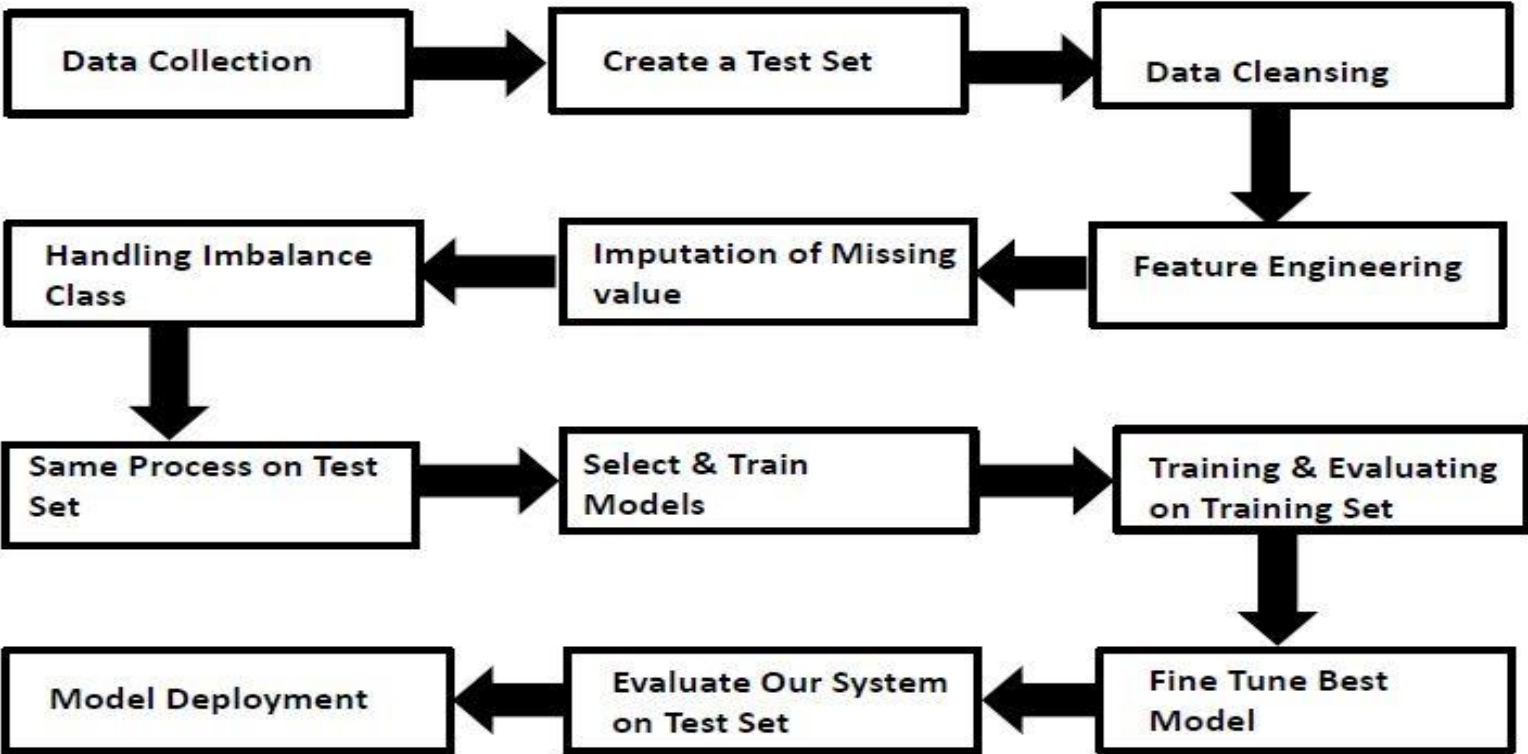
DATASET



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 Simple Imputer)
- 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, SVC , Decision Tree Classifier, KNeighbour are tested.
- Random Forest , Decision Tree and K Neighbour all were given better results. Random Forest was chosen for final model training and testing.
- Model performance evaluated based on accuracy, confusion matrix, classification report.

Random Forest Classifier Model

INTRODUCTION

Random Forest Classifier is an ensemble machine learning algorithm that combines multiple decision trees to make accurate and robust predictions.

The Random Forest Classifier is an ensemble machine learning technique that combines multiple decision trees to enhance to classification accuracy. Each tree contributes its prediction, and the final outcome is determined by majority voting. Known for its robustness and ability to handle various data types, it's widely used in tasks like image classification and fraud detection while providing insights into feature importance.

Reason to use Random Forest Classifier model:

- It has high execution speed.
- It gives better model performance.

MODEL PREDICTION RESULTS ON TEST DATASET

Classification Report

Classification Report for Random Forest				
	precision	recall	f1-score	support
0.0	1.00	0.99	0.99	614
1.0	0.98	1.00	0.99	591
accuracy			0.99	1205
macro avg	0.99	0.99	0.99	1205
weighted avg	0.99	0.99	0.99	1205

[[33	0	0	0]
[2	5	7	0]
[0	0	1068	0]
[0	0	0	1156]]

DATABASE CONNECTION & DEPLOYMENT

Database Connection

- MongoDB Database used for this project.

Model Deployment

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

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 Achieve Folder and a list of these files has been shared with the client and we removed the bad data folder.

Q5) How logs are managed?

We are using logs as per the steps that we follow in training and prediction like model training log and prediction log etc, all the logs are stored under logs 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 ,SVC and Desicion Tree models are trained on all clusters and based on performance, for each cluster different model is saved.

Q 8) How Prediction was done?

- The testing files are shared by the client .We Perform the same life cycle till the data transformation .
- Then model is loaded and perform prediction. In the end we get the accumulated data of predictions.

Q 9) What are the different stages of deployment?

- After model training and model building , We created required files for deployment.
- Finally deployed our model over a cloud platform AWS.

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

- For this project I have made a UI for bulk predictions
- The UI is Very simple and easy to use.
- Client just need to upload csv file then in couple of seconds predictions will be made.

THANK YOU