



**Machine Learning (CS-324)**  
**OEL**

**“Disease Prediction from Symptoms”**

***Submitted to: Miss Mahnoor Malik***

***Prepared by***

<b>GROUP MEMBERS</b>	<b>ROLL NO.</b>
Rooman	CS-21008
Syed Maarij	CS-21022
Usama Khalid	CS-21028

TABLE OF CONTENTS

1. Introduction..... 1

2. Data Collection ..... 1

3. Data Preprocessing..... 1

    Data Cleaning..... 2

4. Exploratory Data Analysis (EDA) ..... 2

5. Feature Engineering ..... 2

6. Model Building ..... 2

7. Model Evaluation..... 2

    Evaluation Metrics ..... 2

        Accuracy: ..... 3

        Confusion Matrix: ..... 3

        Classification Report:..... 4

    Model Comparison..... 4

8. Conclusion ..... 4

9. Limitations and Future Work..... 4

# 1. INTRODUCTION

This report explores the Disease Prediction dataset from Kaggle, performs data preprocessing and exploratory data analysis (EDA), and builds predictive models using various machine learning algorithms. The primary objective is to predict the disease that a person has based on the features available in the dataset. There is an other data 'Disease-Symptom Knowledge Database' which is related to this Disease Prediction dataset, it gives more information. This dataset has 3 columns Disease, Count of Disease Occurrence, Symptom. Disease Prediction dataset used for model training and 'Disease-Symptom Knowledge Database' used for exploratory data analysis (EDA).

## 2. DATA COLLECTION

Applying Knowledge to field of Medical Science and making the task of Physician easy is the main purpose of this Disease Prediction dataset. The disease prediction dataset was downloaded from Kaggle. It contains data for 42 diseases. This dataset has 132 symptoms on which 42 different types of diseases can be predicted. It contains 4920 records. 132 features are;

itching,skin\_rash,nodal\_skin\_eruptions,continuous\_sneezing,shivering,chills,joint\_pain,stomach\_pain,acidity,ulcers\_on\_tongue,muscle\_wasting,vomiting,burning\_micturition,spotting\_urination,fatigue,weight\_gain,anxiety,cold\_hands\_and\_feets,mood\_swings,weight\_loss,restlessness,lethargy,patches\_in\_throat,irregular\_sugar\_level,cough,high\_fever,sunken\_eyes,breathlessness,sweating,dehydration,indigestion,headache,yellowish\_skin,dark\_urine,nausea,loss\_of\_appetite,pain\_behind\_the\_eyes,back\_pain,constipation,abdominal\_pain,diarrhoea,mild\_fever,yellow\_urine,yellowing\_of\_eyes,acute\_liver\_failure,fluid\_overload,swelling\_of\_stomach,swelled\_lymph\_nodes,malaise,blurred\_and\_distorted\_vision,phlegm,throat\_irritation,redness\_of\_eyes,sinus\_pressure,runny\_nose,congestion,chest\_pain,weakness\_in\_limbs,fast\_heart\_rate,pain\_during\_bowel\_movements,pain\_in\_anal\_region,bloody\_stool,irritation\_in\_anus,neck\_pain,dizziness,cramps,bruising,obesity,swollen\_legs,swollen\_blood\_vessels,puffy\_face\_and\_eyes,enlarged\_thyroid,brittle\_nails,swollen\_extremeties,excessive\_hunger,extra\_marital\_contacts,drying\_and\_tingling\_lips,slurred\_speech,knee\_pain,hip\_joint\_pain,muscle\_weakness,stiff\_neck,swelling\_joints,movement\_stiffness,spinning\_movements,loss\_of\_balance,unsteadiness,weakness\_of\_one\_body\_side,loss\_of\_smell,bladder\_discomfort,foul\_smell\_of\_urine,continuous\_feel\_of\_urine,passage\_of\_gases,internal\_itching,toxic\_look\_(typhos),depression,irritability,muscle\_pain,altered\_sensorium,red\_spots\_over\_body,abnormal\_menstruation,dischromic\_patches,watering\_from\_eyes,increased\_appetite,polyuria,family\_history,mucoid\_sputum,rusty\_sputum,lack\_of\_concentration,visual\_disturbances,receiving\_blood\_transfusion,receiving\_unsterile\_injections,coma,stomach\_bleeding,distention\_of\_abdomen,history\_of\_alcohol\_consumption,fluid\_overload,blood\_in\_sputum,prominent\_veins\_on\_calf,palpitations,painful\_walking,pus\_filled\_pimples,blackheads,scurring,skin\_peeling,silver\_like\_dusting,small\_dents\_in\_nails,inflammatory\_nails,blisters,red\_sore\_around\_nose,yellow\_crust\_ooze. The outcome variable is 'prognosis'. The diseases are 'Fungal infection', 'Allergy', 'GERD', 'Chronic cholestasis', 'Drug Reaction', 'Peptic ulcer disease', 'AIDS', 'Diabetes', 'Gastroenteritis', 'Bronchial Asthma', 'Hypertension', 'Migraine', 'Cervical spondylosis', 'Paralysis (brain hemorrhage)', 'Jaundice', 'Malaria', 'Chicken pox', 'Dengue', 'Typhoid', 'hepatitis A', 'Hepatitis B', 'Hepatitis C', 'Hepatitis D', 'Hepatitis E', 'Alcoholic hepatitis', 'Tuberculosis', 'Common Cold', 'Pneumonia', 'Dimorphic hemmorhoids(piles)', 'Heart attack', 'Varicose veins', 'Hypothyroidism', 'Hyperthyroidism', 'Hypoglycemia', 'Osteoarthritis', 'Arthritis', '(vertigo) Parosymal Positional Vertigo', 'Acne', 'Urinary tract infection', 'Psoriasis', 'Impetigo'

## 3. DATA PREPROCESSING

Disease Prediction dataset doesn't need preprocessing but 'Disease-Symptom Knowledge Database' need preprocessing in order to make it easy to visualize.

## Data Cleaning

- **Missing Values:** The Disease and Count of nDisease Occurrence columns had missing values. ffill function was used.
- Remove UMLS codes for diseases and symptoms.
- After doing some processing get the cleaned data that is generated from raw data.

## 4. EXPLORATORY DATA ANALYSIS (EDA)

Plotly, matplotlib and seaborn was used in order to explore data and getting insights from it.

1-Among 42 diseases hypertensive disease has the largest number of occurrences (8.66%).

2- Among 42 diseases decubitus ulcer has the smallest number of occurrences (0.108%).

3- Subplots shows the symptoms responsible for the disease and also shows the count of these symptoms. For example the Gerd disease have symptoms stomach pain, acidity, ulcers on tongue, vomiting, cough and chest pain, among these symptoms stomach pain, cough and chest pain are major symptoms of Gerd disease.

4- Symptoms and associated sub plots shows the all those diseases that have this specific symptom and also shows how much is the importance of that symptom in the disease. For example migraine and gerd have symptom acidity and acidity is the major symptom of migraine disease.

5- Count of each symptom shows fatigue is the major symptom among 42 diseases.

## 5. FEATURE ENGINEERING

Created 3 feature Disease, Count of Disease Occurrence, Symptom by doing processing on 'Disease-Symptom Knowledge Database'.

## 6. MODEL BUILDING

### Machine Learning Algorithms

- **Decision Trees:** Implemented using scikit-learn and from scratch.
- **Random Forest:** Implemented using scikit-learn and from scratch.

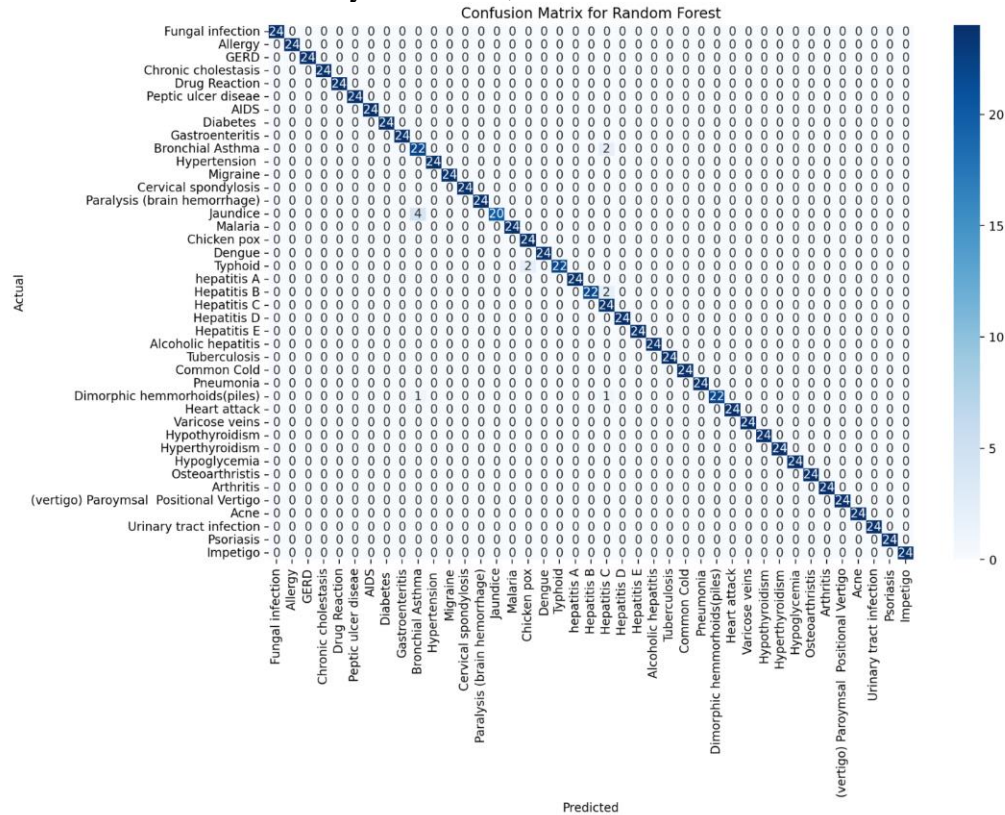
## 7. MODEL EVALUATION

### Evaluation Metrics

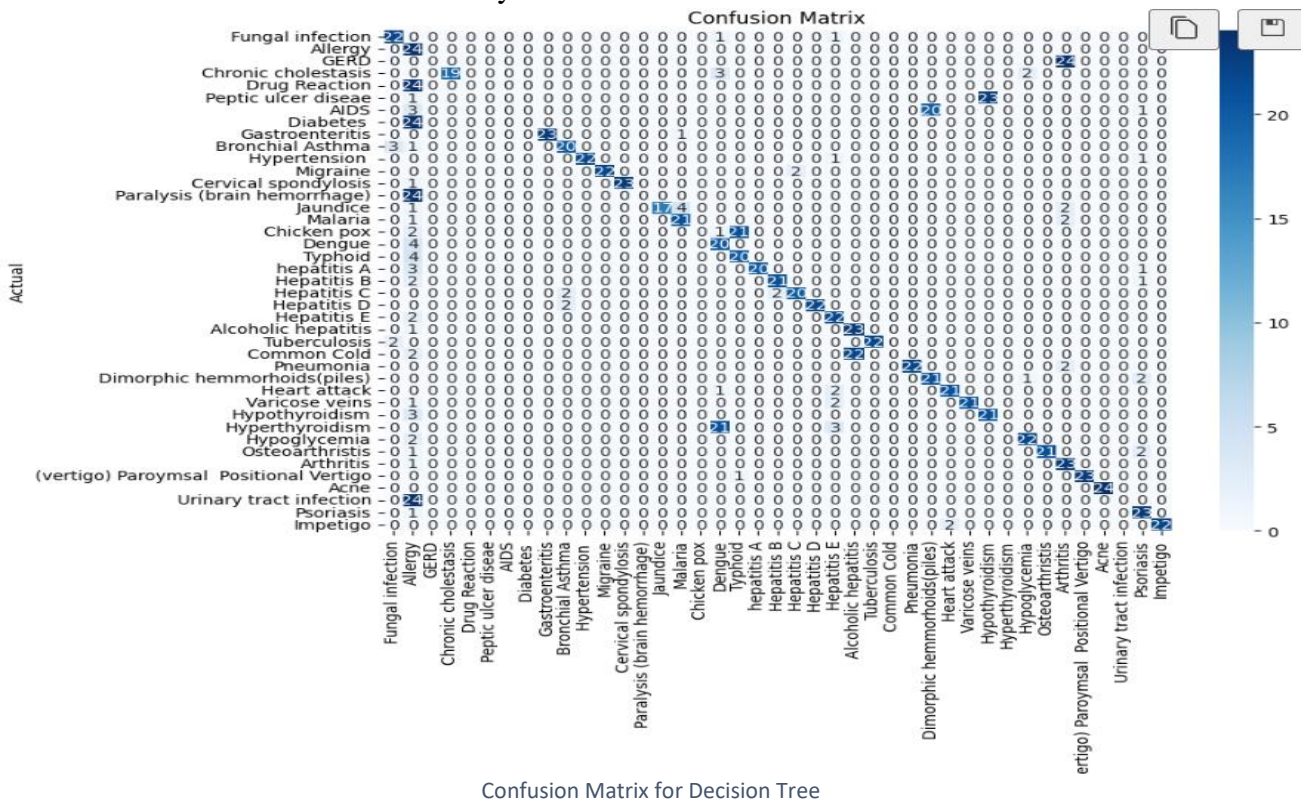
- **Accuracy:** Measured the overall accuracy of each model.
- **Confusion Matrix:** Evaluated the performance by understanding the true positives, false positives, true negatives, and false negatives.
- **Classification Report:** Provided precision, recall, and F1-score for each class.

## Model Comparison

- Random Forest: Achieved an accuracy of 97.62%,



- Decision Tree: Achieved an accuracy of 67.78%



## • CLASSIFICATION REPORT

Classification Report:

precision recall f1-score support

Fungal infection	0.81	0.92	0.86	24
Allergy	0.15	1.00	0.27	24
GERD	0.00	0.00	0.00	24
Chronic cholestasis	1.00	0.79	0.88	24
Drug Reaction	0.00	0.00	0.00	24
Peptic ulcer disease	0.00	0.00	0.00	24
AIDS	0.00	0.00	0.00	24
Diabetes	0.00	0.00	0.00	24
Gastroenteritis	1.00	0.96	0.98	24
Bronchial Asthma	0.83	0.83	0.83	24
Hypertension	1.00	0.92	0.96	24
Migraine	1.00	0.92	0.96	24
Cervical spondylosis	1.00	0.96	0.98	24
Paralysis (brain hemorrhage)	0.00	0.00	0.00	24
Jaundice	1.00	0.71	0.83	24
Malaria	0.81	0.88	0.84	24
Chicken pox	0.00	0.00	0.00	24
Dengue	0.43	0.83	0.56	24
Typhoid	0.48	0.83	0.61	24
hepatitis A	1.00	0.83	0.91	24
Hepatitis B	0.91	0.88	0.89	24
Hepatitis C	0.91	0.83	0.87	24
Hepatitis D	1.00	0.92	0.96	24
Hepatitis E	0.71	0.92	0.80	24
Alcoholic hepatitis	0.51	0.96	0.67	24
Tuberculosis	1.00	0.92	0.96	24
Common Cold	0.00	0.00	0.00	24
Pneumonia	1.00	0.92	0.96	24
Dimorphic hemmorhoids(piles)	0.51	0.88	0.65	24
Heart attack	0.91	0.88	0.89	24
Varicose veins	1.00	0.88	0.93	24
Hypothyroidism	0.48	0.88	0.62	24
Hyperthyroidism	0.00	0.00	0.00	24
Hypoglycemia	0.88	0.92	0.90	24
Osteoarthritis	1.00	0.88	0.93	24
Arthritis	0.43	0.96	0.60	24
(vertigo) Paroymsal Positional Vertigo	1.00	0.96	0.98	24
Acne	1.00	1.00	1.00	24
Urinary tract infection	0.00	0.00	0.00	24
Psoriasis	0.74	0.96	0.84	24
Impetigo	1.00	0.92	0.96	24
accuracy			0.68	984
macro avg	0.62	0.68	0.63	984
weighted avg	0.62	0.68	0.63	984

Classification Report on Test Data for Random Forest:

precision recall f1-score support

0	1.00	1.00	1.00	1
1	1.00	1.00	1.00	1
2	1.00	1.00	1.00	1
3	1.00	1.00	1.00	1
4	1.00	1.00	1.00	1
5	1.00	1.00	1.00	1
6	1.00	1.00	1.00	1
7	1.00	1.00	1.00	1
8	0.50	1.00	0.67	1
9	1.00	1.00	1.00	1
10	1.00	1.00	1.00	1
11	1.00	1.00	1.00	1
12	1.00	1.00	1.00	1
13	1.00	1.00	1.00	1
14	1.00	1.00	1.00	1
15	1.00	0.50	0.67	2
16	1.00	1.00	1.00	1
17	1.00	1.00	1.00	1
18	1.00	1.00	1.00	1
19	1.00	1.00	1.00	1
20	1.00	1.00	1.00	1
21	1.00	1.00	1.00	1
...				
accuracy			0.98	42
macro avg	0.99	0.99	0.98	42
weighted avg	0.99	0.98	0.98	42

## 8. CONCLUSION

For implementing disease prediction model, use random forest because it has high accuracy 98%. The insights come out from EDA are hypertensive disease has the largest number of occurrences, decubitus ulcer has the smallest number of occurrences and fatigue is the major symptom among 42 diseases.

## 9. LIMITATIONS AND FUTURE WORK

### Limitations:

1. Dataset Size: Limited to 4920 records and 42 diseases, which may not capture full variability.
2. Symptom Overlap: Common symptoms among diseases can lead to misclassification.
3. Symptom Details: Severity and duration of symptoms are not considered.
4. Static Data: The model does not adapt to new data over time.
5. Data Quality: Potential remaining noise or inaccuracies in the dataset.

### Future Work:

1. Larger Datasets: Incorporate more comprehensive datasets with more records and diseases.
2. Temporal Analysis: Include symptom progression and duration.
3. Real-time Integration: Enable real-time data updates.
4. Advanced Features: Explore sophisticated feature engineering techniques.
5. Model Interpretability: Improve model explanations for healthcare professionals.
6. Multi-modal Data: Integrate genetic, imaging, and patient history data.

