

## **CHAPTER-01**

### **COVID-19 Diagnosis Prediction Using Machine Learning**

#### **1.1 INTRODUCTION**

The COVID-19 Diagnosis Prediction Project leverages machine learning to predict COVID-19 infection outcomes based on patient symptoms and demographic data. This project is rooted in the healthcare field, focusing on the early detection of infectious diseases and efficient resource management. The COVID-19 pandemic posed significant challenges, especially in resource allocation and containment. The ability to predict an individual's infection status can streamline resource planning, improve patient outcomes, and reduce the healthcare system's burden.

The project employs technologies such as Python for data analysis, and various machine learning algorithms like Logistic Regression, Decision Trees, and XGBoost. These tools facilitate data cleaning, feature selection, and model training to develop accurate predictive models. Special terms like Exploratory Data Analysis (EDA), hyperparameter tuning, and feature engineering are crucial components of the methodology.

By implementing this project, healthcare providers can enhance screening efficiency, optimize resource utilization, and potentially extend the approach to predict outcomes for other diseases in the future.

## 1.2 RATIONALE

Our project is crucial for early COVID-19 detection, efficient healthcare resource allocation, and informed public health decision-making, ultimately improving overall well-being.

**Q1.** Why is your proposal important in today's world? How predicting a disease accurately can improve medical treatment?

**ANS :** Accurately predicting diseases using machine learning, particularly for COVID-19, is highly significant in today's world. Early detection based on symptoms enables healthcare facilities to isolate and treat patients promptly, which is crucial for preventing the spread of the disease.

**Q2.** How is it going to impact the medical field when it comes to effective screening and reducing health care burden.

**ANS:** Impact on Medical Facilities: Currently, medical facilities use tests such as the Rapid Antigen Test (RAT) to detect COVID-19. Our machine learning model leverages patient data to accurately determine COVID-19 infection, significantly reducing the burden of testing and screening on healthcare facilities.

**Q3.** If any, what is the gap in the knowledge or how your proposed method can be helpful if required in future for any other disease.

**ANS:** Future Applications: Poor-quality or missing data can pose challenges for these models. However, if our model proves effective for COVID-19, similar approaches could be applied to other diseases in the future, providing a valuable tool for medical facilities.

### 1.3 PROBLEM STATEMENT

The COVID-19 pandemic created a significant global health crisis, challenging healthcare systems worldwide. The rapid spread of the virus placed an immense burden on medical resources, including hospital beds, ventilators, and personal protective equipment (PPE). One of the major issues faced by healthcare providers was the inability to predict resource requirements in advance, leading to shortages and delays in treatment. This situation emphasized the urgent need for predictive tools that can assist in resource allocation and patient management.

Machine learning-based models can play a pivotal role in addressing this challenge. By analyzing symptoms, demographic data, and patient history, these models can predict whether a patient is likely to test positive for COVID-19. Early identification of high-risk patients allows for better preparation and resource distribution. Moreover, such models can guide healthcare providers in prioritizing patients based on the severity of symptoms and anticipated resource needs. In addition to enhancing patient outcomes, predictive tools reduce the burden on testing facilities and improve the overall efficiency of healthcare systems. They enable targeted testing, minimizing unnecessary tests and focusing resources on individuals most likely to benefit. This streamlined approach is critical in managing the spread of the virus and optimizing healthcare operations .

Ultimately, the application of machine learning in disease prediction extends beyond COVID-19. It highlights the potential for similar models to be used in predicting and managing other infectious diseases, contributing to a more proactive and data-driven healthcare system.

## 1.4 OBJECTIVE

- To efficiently predict and manage COVID-19 cases.
- To analyze patient symptoms and demographics using machine learning.
- To enable early detection, resource allocation, and faster healthcare responses.
- To provide insights that can improve handling of future health crises.
- To leverage technology for meaningful advancements in public health.

## **1.5 PROPOSED METHODOLOGY**

The proposed methodology for the COVID-19 Diagnosis Prediction Project involves data-driven approaches and machine learning techniques to achieve accurate predictions and resource optimization.

### **1.5.1 Research Type**

This project follows an applied research approach, focusing on the practical application of machine learning to predict COVID-19 outcomes using real-world data.

#### **Data Collection and Analysis**

Data Source: Patient data, including symptoms and demographic information, is collected from publicly available datasets.

Tools: Python is used for data processing, analysis, and machine learning model development. MySQL manages the data storage and retrieval.

Methods: Exploratory Data Analysis (EDA) is performed to clean, visualize, and understand the data, followed by feature engineering to enhance model performance.

#### **Steps to Achieve the Objective**

Step 1: Data Cleaning: Handle missing values, remove duplicates, and prepare the dataset.

Step 2: Feature Engineering: Select and transform features relevant to the prediction task.

Step 3: Model Training: Train multiple machine learning models (Logistic Regression, Decision Trees, Random Forest, XGBoost).

Step 4: Hyperparameter Tuning: Optimize the models for accuracy and performance.

Step 5: Model Evaluation: Use metrics such as accuracy, precision, and recall to evaluate the models' performance.

This systematic approach ensures robust and accurate predictions, enabling better healthcare planning and decision-making.

### 1.5.2 Justification for Train-Test Splitting

Purpose:

- To evaluate the model's performance on unseen data.
- To ensure the model generalizes well and does not overfit to the training data.

Key Advantages:

- Simplicity: Straightforward to implement and interpret.
- Efficiency: Requires less computation compared to techniques like cross-validation.
- Control Over Testing Size: Allows you to adjust the test set size (test\_size parameter).

Consistency:

- By setting random\_state, you ensure reproducibility in results.

### 1.5.3 Why Not Use Other Splitting Techniques?

**K-Fold Cross-Validation:**

- Splits the data into multiple folds for training and validation.
- Advantage: More robust performance evaluation.
- Disadvantage: Higher computational cost, which may not be necessary for this task.

**Leave-One-Out Cross-Validation:**

- Each observation acts as a test set once.
- Advantage: Utilizes the maximum amount of data for training.
- Disadvantage: Computationally expensive and may not be practical for large datasets.

**Stratified Splitting:**

- Ensures class proportions in training and test sets match the original dataset.
- While this is beneficial for imbalanced datasets, it adds complexity and may not be crucial if the dataset is balanced.

In our case, train\_test\_split is a suitable choice as:

- It strikes a balance between computational efficiency and ensuring model generalization.
- The dataset size is large enough to allocate 20% for testing without sacrificing training data size.

### 1.5.4 Justification for Selected Algorithms

#### **Logistic Regression:**

- Simple and efficient for binary classification.
- Serves as a baseline with low computational cost and interpretability.

#### **Decision Tree:**

- Handles both categorical and numerical data.
- Captures non-linear relationships and provides an interpretable model.

#### **Random Forest:**

- Reduces overfitting of decision trees through ensemble learning. □ Robust to noise and provides feature importance metrics.

#### **XG Boost Classifier:**

- Highly accurate with gradient boosting and regularization.
- Efficient handling of missing data and strong feature importance insights.

### 1.5.5 Why Not Other Algorithms?

- KNN: Inefficient for large datasets and lacks interpretability.
- SVM: Computationally expensive and harder to tune for large datasets. □ Naive Bayes: Assumes feature independence, which is unrealistic here.
- Neural Networks: Overkill for structured data; requires extensive tuning.

These selected algorithms balance complexity, accuracy, and computational efficiency.

### 1.5.6 Step-by-Step Process for COVID-19 Data Analysis and Model Training

#### **1.Import Libraries:**

- Libraries like numpy, pandas, matplotlib, and seaborn are imported for data manipulation, visualization, and statistical analysis.
- Machine learning libraries like scikit-learn and xgboost are also imported for model training and evaluation.

## 2. Load the Dataset:

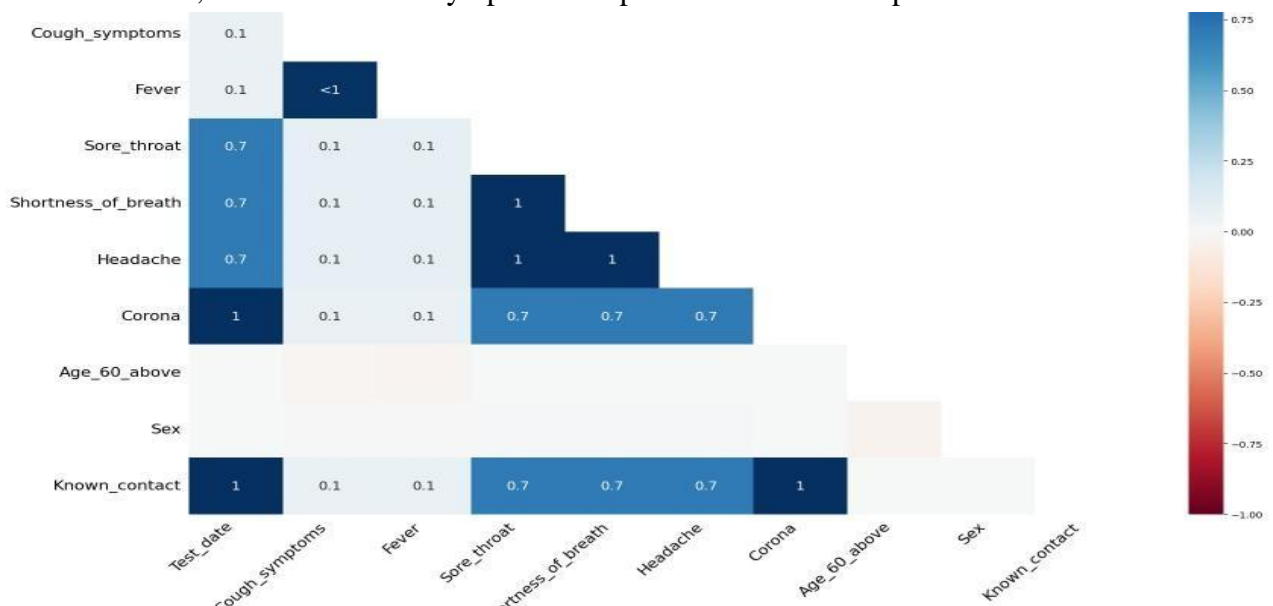
- The dataset is read from a CSV file (corona\_tested\_006.csv) into a pandas DataFrame for analysis.
- A copy of the dataset is created to preserve the original data.

## 3. Exploratory Data Analysis (EDA):

- The dataset shape, column details, and unique values are inspected. ☐ Missing values are analyzed, and their percentage is calculated.
  - To ensure data quality by identifying incomplete data, guide cleaning decisions (imputation or removal), confirm model compatibility to prevent errors, and maintain consistent inputs for improved model performance.
- Visualizations like heatmaps are used to study correlations and missing data patterns.

### Correlation Heatmap Analysis

- A heatmap was generated using `sns.heatmap(df_train_val.corr(), annot=True)`.
- It shows the pairwise correlation between numerical features in the dataset.
- High correlation values (closer to 1 or -1) suggest a strong relationship between features. For instance, if Cough\_symptoms has a strong positive correlation with Corona, it indicates these symptoms are prevalent in COVID-positive cases.



**Figure 1.1:** Correlation Heatmap Analysis

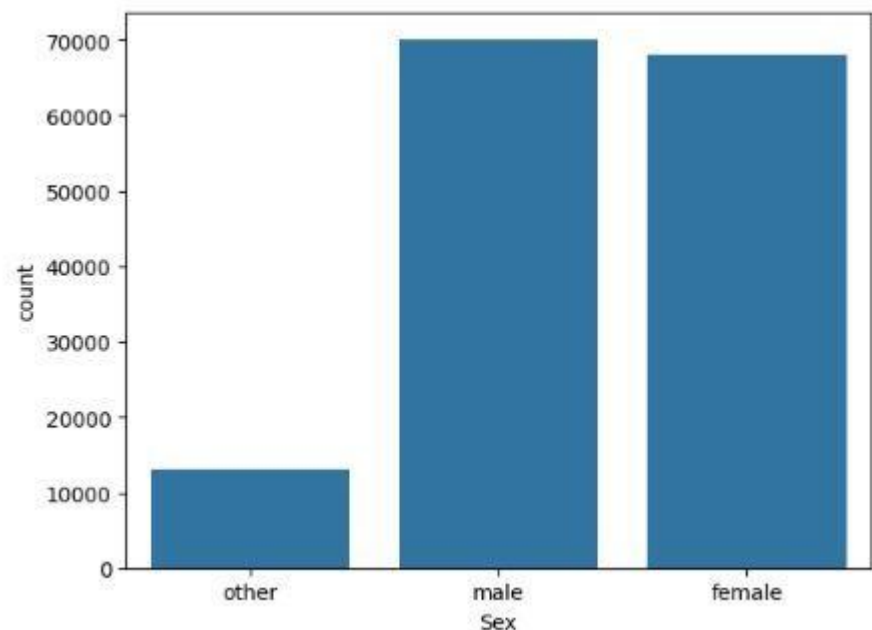


#### 4. Data Cleaning:

□ Null values are either imputed or dropped, depending on their significance. □ Some columns, such as Sex, are filled with default values like 'other'.

##### Gender Distribution in the Dataset

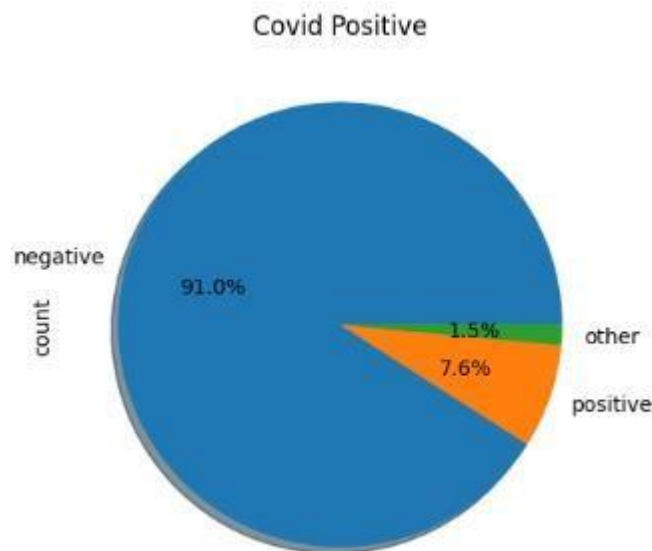
- Description: The countplot depicts the distribution of individuals across three gender categories: Male, Female, and Other.
- Observations:
  - The Male and Female categories have a similar number of individuals, with both close to 70,000.
  - The Other category represents a significantly smaller portion of the population, with approximately 10,000 individuals.
  - Insights:
    - This balanced representation of males and females ensures that gender-related analysis for COVID diagnosis will not be biased.
    - The smaller sample size for the Other category might make it harder to derive conclusive insights for this group.



**Figure 1.2:** Gender Distribution in the Dataset

**Pie Chart:**

- `covid['Corona'].value_counts().plot.pie()` was used to show the distribution of COVID-positive, negative, and other cases.
- This chart provides a proportional overview of the target variable, helping in understanding dataset balance.



**Figure 1.3:** Pie Chart for COVID-19 case analysis

**Feature Distribution via Count Plots****Distribution of COVID Cases by Cough Symptoms**

- Purpose: To understand how the presence of cough symptoms correlates with COVID test results.
- Insights:
  - Bars for True in Cough\_symptoms likely show a higher count in Corona: positive compared to negative.
  - Helps in identifying whether Cough\_symptoms is a significant indicator of COVID positivity.

**Distribution of COVID Cases by Fever**

- Purpose: To analyze the impact of fever on COVID test outcomes.
- Insights:
  - Bars for True in Fever under Corona: positive indicate its prevalence among positive cases.

- Highlights fever as a potential symptom for diagnosis.

#### Distribution of COVID Cases by Sore Throat

- Purpose: To determine the relationship between sore throat and COVID results.
- Insights:
  - Counts for True in Sore\_throat might show a pattern where this symptom is less frequent in COVID-positive cases compared to others.

#### Distribution of COVID Cases by Shortness of Breath

- Purpose: To explore whether shortness of breath is strongly associated with positive COVID cases.
- Insights:
  - Bars for True in Shortness\_of\_breath could reveal its diagnostic significance in identifying severe cases.

#### Distribution of COVID Cases by Headache

- Purpose: To assess how headaches correlate with COVID test outcomes.
- Insights:
  - Bars for True in Headache under Corona: positive may highlight whether headaches are common in COVID-positive patients.

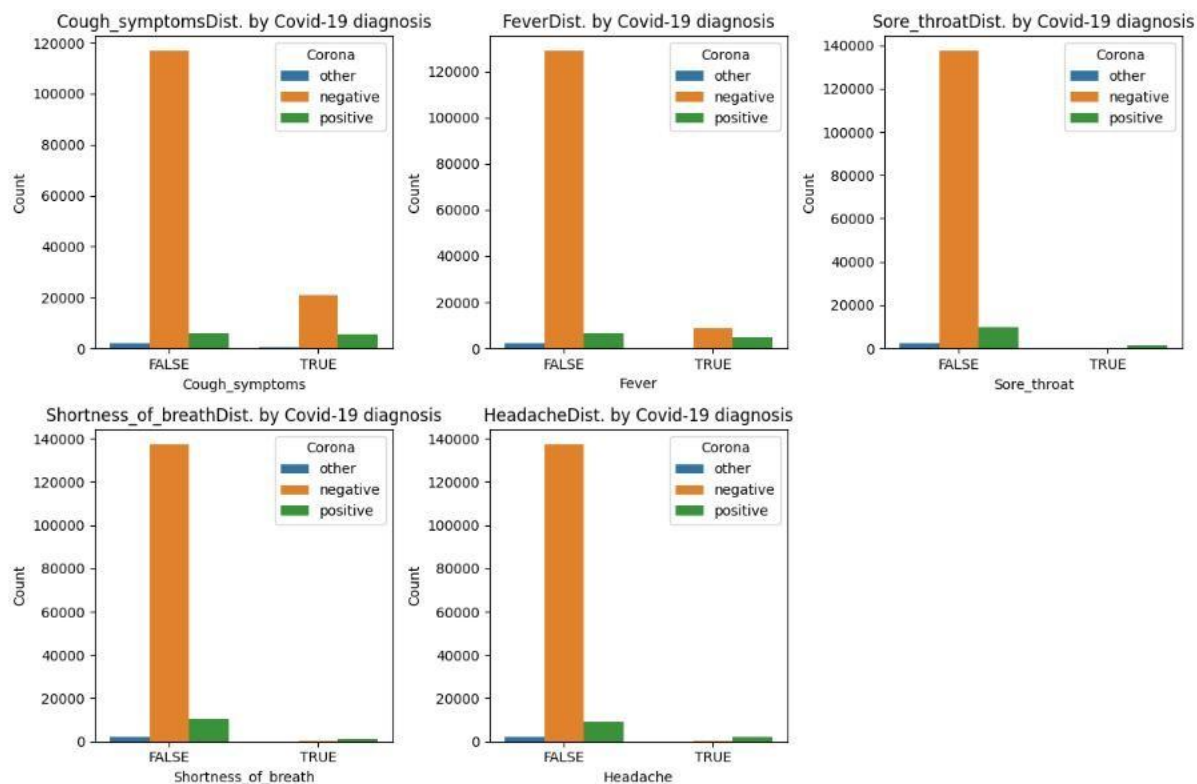
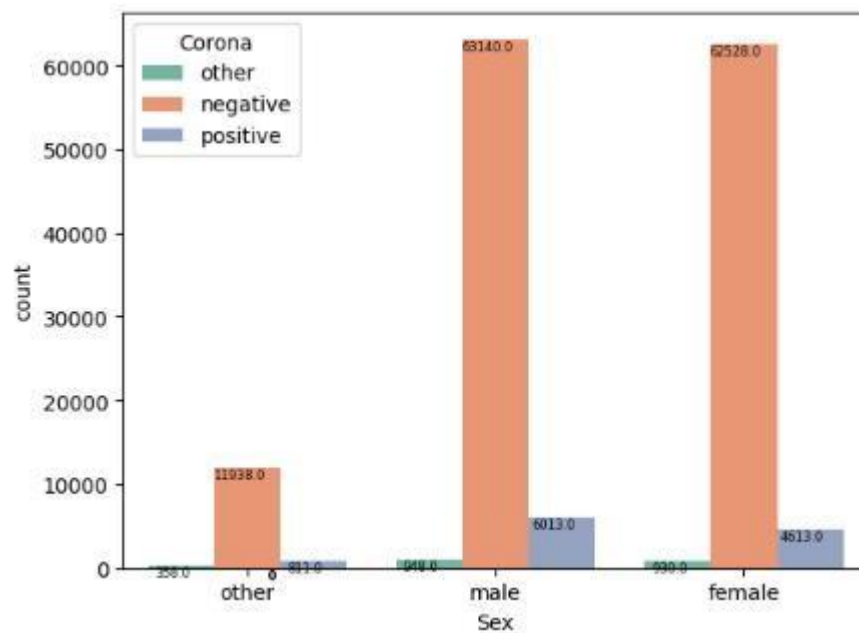


Figure 1.4: Feature Distribution via Count Plots

## Gender and Fever Distribution by COVID Status

### 1. Gender Distribution by COVID Status

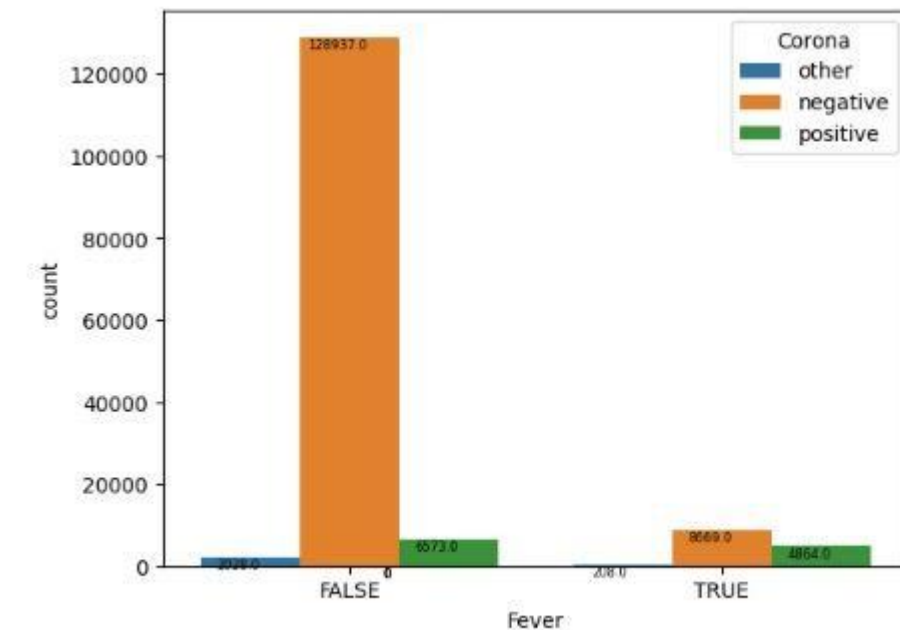
- Description: The top countplot illustrates the distribution of COVID statuses (Positive, Negative, Other) across three gender categories (Male, Female, and Other).
- Observations:
  - Both males and females have a significantly higher count of COVID-negative cases compared to positive and other statuses.
  - COVID-positive cases are relatively small across all genders but slightly higher in males and females compared to the "Other" gender category. □ Insights:
  - The imbalance between negative and positive cases reflects the overall dataset composition.
  - This distribution shows no significant gender-based difference in COVID positivity rates.



**Figure 1.5:** Gender Distribution by COVID Status

## 2. Fever Distribution by COVID Status

- Description: The bottom countplot shows the relationship between fever presence (True or False) and COVID statuses (Positive, Negative, Other). □ Observations:
  - A majority of COVID-negative cases are associated with "Fever: False."
    - COVID-positive cases are more evenly distributed between "Fever: True" and "Fever: False," but their absolute count remains low compared to negatives.
- Insights:
  - Fever alone is not a definitive indicator of COVID positivity but might have some predictive value in combination with other symptoms.
  - The higher prevalence of "Fever: False" across all COVID statuses suggests fever is not a universal symptom.



**Figure 1.6:**Fever Distribution by COVID Status

### 5. Feature Engineering:

- Columns like Test\_date are converted to appropriate formats (e.g., datetime).
- Boolean features (e.g., Cough\_symptoms) are encoded to TRUE/FALSE strings for consistency.

### 6. Data Splitting:

- The dataset is split into training-validation (df\_train\_val) and testing (df\_test) sets based on the date range.
- Unnecessary columns like Test\_date, Ind\_ID, and Sex are dropped.

### 7. Label Encoding:

- ❑ Categorical variables (e.g., Cough\_symptoms, Fever) are converted into numerical labels using LabelEncoder for model compatibility.

### 8. Feature Scaling:

- ❑ The StandardScaler is applied to normalize the training data to ensure uniform scaling of features.

## 9. Feature Selection:

- The Chi-Square test is performed to evaluate feature importance.
- A heatmap is used to visualize feature correlations.

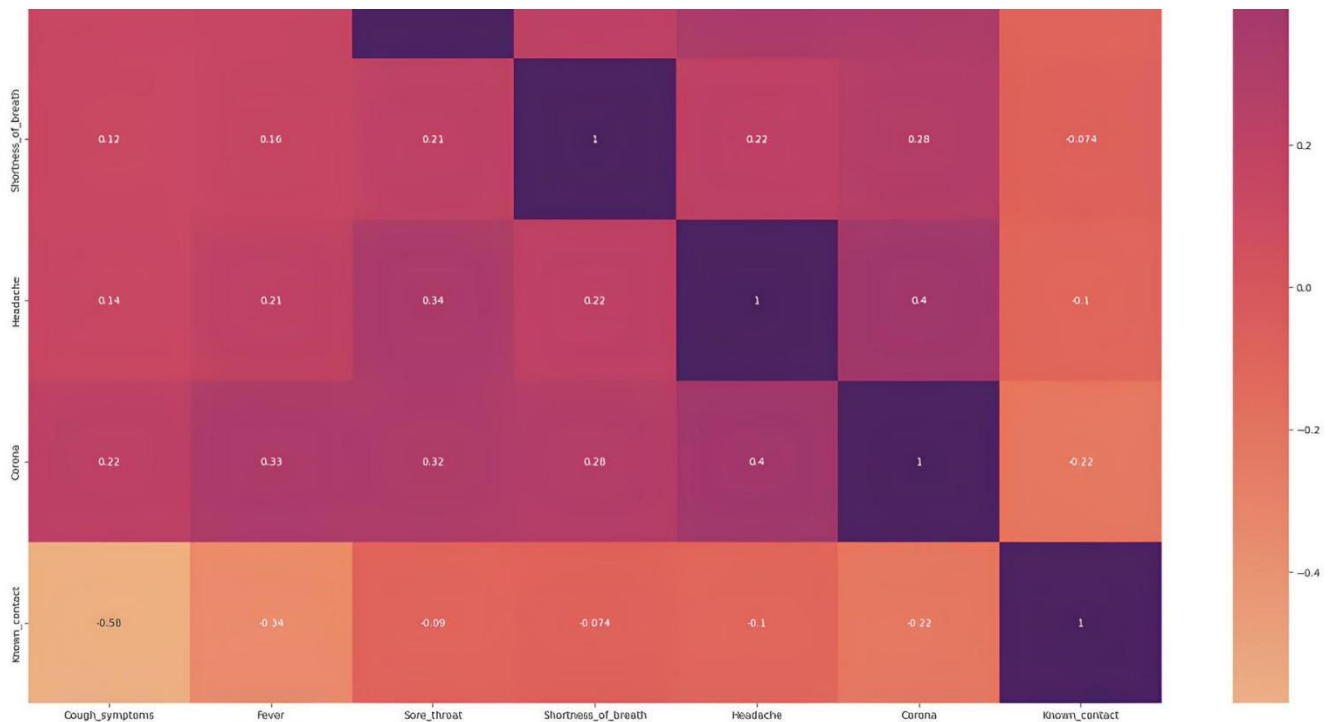


Figure 1.7: Correlation Heatmap Analysis

## 10. Model Training and Validation:

- Logistic Regression: Trained on the data to act as a baseline model.
- Decision Tree: Built to capture non-linear relationships.
- Random Forest: Used for ensemble learning to improve accuracy and robustness.
- XGBoost Classifier: Applied for advanced gradient boosting to achieve high accuracy.

## 11. Model Testing:

- Each trained model is evaluated on the testing dataset for accuracy and other metrics.

## 12. Evaluation Metrics:

- Metrics such as accuracy,  $R^2$  score, and RMSE are calculated for all models to compare their performance.

## 1.6 EXPECTED OUTCOME

The expected outcome of this project is the development of a machine learning-based predictive model capable of accurately diagnosing COVID-19 based on symptoms and demographic data. The model will enable early detection of potential cases, thereby improving the efficiency of healthcare resource allocation and patient management. By achieving high accuracy and reliability, the project aims to support clinical decision-making, reduce the burden on testing facilities, and enhance overall healthcare outcomes. This predictive tool can also serve as a foundation for similar applications in diagnosing other infectious diseases in the future, contributing to proactive healthcare strategies.

### Model Comparison

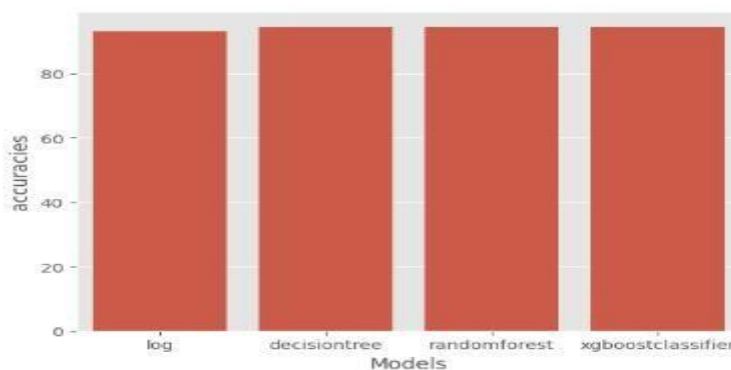
Data:

The table shows four models and their corresponding accuracies:

- Logistic Regression (log): 93.14%
- Decision Tree (decisiontree): 94.46%
- Random Forest (randomforest): 94.46%
- XGBoost Classifier (xgboostclassifier): 94.46%

Visualization:

A bar plot is used to visually compare the accuracies of these models. The x-axis represents the models, and the y-axis shows the accuracy percentages. All models perform similarly, with Decision Tree, Random Forest, and XGBoost yielding identical accuracy of 94.46%, slightly outperforming Logistic Regression.



**Figure 1.8:** Model Comparison



## CONCLUSION

This investigation provides an insightful comparison of machine learning models, highlighting the similar performance of Decision Tree, Random Forest, and XGBoost Classifier, each with an accuracy of 94.46%, and Logistic Regression, which closely followed with 93.14%. The comparable outcomes of the tree-based models suggest their effectiveness in handling the complexities of the dataset, possibly due to inherent capabilities like feature importance evaluation and handling non-linearity.

However, the study's preliminary nature—conducted without cross-validation or hyperparameter tuning—suggests that the reported results may not fully capture the models' potential. Future work should focus on enhancing the evaluation process by employing rigorous cross-validation to reduce the risk of overfitting and applying hyperparameter tuning to maximize performance. Additionally, incorporating other metrics, such as area under the curve (AUC) and confusion matrix analysis, will provide a comprehensive understanding of the models' strengths and weaknesses.

Ultimately, this study illustrates the importance of iterative refinement in machine learning workflows. By leveraging tools like Python's pandas for data handling and seaborn for visualization, researchers can efficiently compare models and make informed decisions for model selection. The demonstrated process serves as a valuable framework for similar classification tasks, particularly in healthcare and related domains.

## REFERENCES

### Datasets

COVID-19 Data Repository:

Johns Hopkins University provides global COVID-19 data, including confirmed cases, deaths, and recoveries.

Access link: [COVID-19 Data Repository by JHU](#)

UCI Machine Learning Repository:

COVID-19 symptom datasets, including demographic details and medical history.

Access link: [UCI Repository - COVID-19](#)

### Research Papers

Y. Zoabi et al. "Machine learning-based prediction of COVID-19 diagnosis," discusses using binary features for prediction.

M. Pourhomayoun et al., "Predicting mortality risk in patients with COVID-19," presents an AI model for healthcare prioritization.

C. Iwendi et al., "COVID-19 Patient Health Prediction Using Boosted Algorithms," explores ML in healthcare systems.

### Machine Learning Documentation

Scikit-Learn Documentation: Comprehensive guides on data preprocessing, model training, and hyperparameter tuning.

Access link: [Scikit-Learn Documentation](#)

TensorFlow Documentation: Details on deep learning models and their application to COVID-19 datasets.

Access link: [TensorFlow Documentation](#)

These resources will help in data analysis, model building, and evaluating the performance of your machine learning models.

