

User Documentation for Analysing the Performance of Machine Learning Algorithms For COVID-19 Symptoms Checker

Overview

The project aims to evaluate various machine learning algorithms to find the most accurate and reliable methods for predicting COVID-19 based on symptoms. By comparing algorithms like decision trees, random forests, neural networks, and logistic regression, the goal is to improve the effectiveness of COVID-19 symptom checkers. This involves analyzing large datasets to identify patterns and develop tools for early detection and better management of COVID-19. The research seeks to determine which machine learning technique works best for symptom checking, helping doctors quickly identify COVID-19 cases, which could enhance patient outcomes and reduce the virus's spread.

1. Prerequisites

To get started, make sure you have the following:

□ Python Version: 3.x

□ Required Libraries:

- pandas for working with data
- scikit-learn for machine learning tasks
- matplotlib and seaborn for creating visualizations

Install all the necessary libraries with this command:

- “pip install pandas scikit-learn matplotlib seaborn”
-

2. Data Collection and Pre-processing

2.1. Data Collection

- The dataset was originally compiled by public health authorities in India. It was collected between from various healthcare facilities, including hospitals, clinics, and testing centres across the country. The data collection process involved gathering clinical information from patients who were evaluated for COVID-19, either because they exhibited symptoms or had been in contact with confirmed cases.

2.2. Data Pre-processing

- In data preprocessing, the main steps involve cleaning and preparing the dataset to ensure it is suitable for machine learning models. This typically includes handling missing data, scaling numerical features to standardize their range, and encoding

categorical variables into numerical values. Techniques like SMOTE may be used to address class imbalances in the dataset, ensuring that the model has a balanced exposure to all classes during training. These preprocessing steps are crucial to improve the model's performance and reliability in making accurate predictions.

3. Model Training and Evaluation

3.1. Linear Regression

- Purpose: Logistic Regression was chosen as a primary model due to its clearness and ease of interpretation. It is very commonly used model for binary classification tasks and provides insights into the relationship between features and the target variable through its coefficients.
- Evaluation: The model's performance is evaluated relationship between COVID-19 symptoms and the likelihood of infection, providing a baseline model to predict the presence of COVID-19 based on clinical and demographic data. It helps assess the linear associations between features and the outcome.

3.2. Random Forest Regression

- Purpose: I selected the Random Forest classifier for its robustness and ability to manage complex, non-linear relationships within the data. Random Forest is an ensemble learning method that constructs number of decision trees and aggregates their predictions to improve correctness and control overfitting.
- Evaluation: Random Forest evaluates complex, non-linear relationships between COVID-19 symptoms and infection likelihood by aggregating predictions from multiple decision trees. It helps identify the most important features and enhances prediction accuracy through its ensemble approach.

3.3 Decision Tree Regression

- Purpose: The Decision Tree classifier was chosen for its interpretability. Decision Trees provide transparent, visible representation of the decision-making process, making it easy to trace how the model arrives at a particular prediction
- Evaluation: The Decision Tree evaluates how individual COVID-19 symptoms contribute to the likelihood of infection by creating a simple, interpretable model that visually maps decision paths. It captures non-linear relationships and provides clear insights into the decision-making process.

3.4 Model Comparison

- Purpose: Compares the predictive power of the Linear Regression and Random Forest and Decision Tree models to identify which model performs better.
-

4. Data Visualization

Here are the types of data visualizations used in this project:

1. Histograms: To visualize the distribution of individual features, such as age or symptom frequency.
2. Scatter Plots: To explore the relationships between different variables, like symptoms and demographic data.
3. Heatmaps: To display correlations between variables, helping identify strong and weak relationships.
4. Confusion Matrix: To visualize the performance of the classification models in terms of true positives, true negatives, false positives, and false negatives.
5. ROC Curve: To illustrate the trade-off between sensitivity and specificity for the different models, providing a measure of model performance.
6. Feature Importance Plot: To rank and visualize the importance of different features in the Random Forest model.

5. Learning Curves

- Purpose: Plots learning curves for both the Linear Regression and Random Forest and Decision tree models to see how performance changes as the training data size increases.
- Insight: Learning curves help identify whether the models are overfitting or underfitting, giving you a deeper understanding of their generalization capabilities.

6. Conclusion

This project demonstrated the effectiveness of using machine learning, particularly the Random Forest model, to predict COVID-19 based on medical symptoms and demographic data. The Random Forest model outperformed others, proving highly reliable for early detection and patient triage. However, future research should focus on incorporating more data, validating the model across diverse populations, and addressing ethical considerations to ensure its practical application in real-world healthcare settings. The findings highlight the potential of machine learning in improving public health responses during pandemics.