# MEDICAL DIAGNOSIS WITH NAIVE BAYES

## A PROJECT REPORT

for

**Introduction to AI(AI101B)**

**Session (2024-25)**

**Submitted by**

**Arun Kumar (202410116100040 )**

**Ayushi Saran Singh(202410116100047 )**

**Submitted in partial fulfilment of the**

**Requirements for the Degree of**

## MASTER OF COMPUTER APPLICATION

**Under the Supervision of**

**Mr. Apoorv Jain**

**Assistant Professor**



**Submitted to**

**Department Of Computer Applications**

**KIET Group of Institutions, Ghaziabad**

**Uttar Pradesh-201206**

**(April - 2025)**

# CERTIFICATE

Certified that Arun Kumar (202410116100040) and Ayushi Saran Singh (202410116100047) have carried out the project work having "**Medical Diagnosis with Naive Bayes**" (INTRODUCTION TO AI) (AI101B)) for **Master of Computer Application** from Dr. A.P.J. Abdul Kalam Technical University (AKTU**)** (formerly UPTU), Lucknow under my supervision. The project report embodies original work, and studies are carried out by the student himself/herself and the contents of the project report do not form the basis for the award of any other degree to the candidate or to anybody else from this or any other University/Institution.

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Date : 19 – 04 - 2025

|  |  |
|---|---|
| **Mr. Apoorv Jain** | **Dr. Akash Rajak** |
| **Assistant professor** | **Dean** |
| **Department of Computer Applications** | **Department of Computer Applications** |
| **KIET Group of Institutions, Ghaziabad** | **KIET Group of Institutions, Ghaziabad** |
| **An Autonomous Institutions** | **An Autonomous Institutions** |

# DECLARATION

I hereby declare that the work presented in this report entitled **"Medic: An AI-Based Medical Diagnosis System using Naive Bayes"** has been carried out by us under the supervision of **Mr. Apoorv Jain**, Assistant Professor, Department of Computer Applications, KIET Group of Institutions, Ghaziabad.

I further declare that:

- This work is original and has not been submitted to any other university or institution for the award of any degree or diploma.

- Due credit has been given to all sources used, including ideas, graphics, codes, and results from previously published works.

- I have used quotation marks wherever direct sentences or statements have been cited, and all non-original work has been properly referenced.

- No part of this project report is plagiarized. In case of any complaint regarding plagiarism or manipulation of data, I shall be held fully responsible and accountable.

**Name**: Arun Kumar , Ayushi Saran Singh
**University Roll No.**: 202410116100040, 202410116100047
**Branch**: Master of Computer Application (MCA)

# ABSTRACT

The project titled *"Medic"* is an AI-based medical diagnosis system developed using the Naive Bayes classification algorithm. Its primary objective is to assist healthcare professionals and individuals in determining the severity of illness based on symptoms and demographic data, thereby enabling timely medical intervention. The system is built on a supervised machine learning model trained on a symptom-based dataset comprising over 63,000 patient records with features such as fever, dry cough, sore throat, tiredness, breathing difficulty, age group, gender, and recent contact with infected individuals.

The motivation behind this project stems from the need for accessible, scalable, and accurate preliminary diagnosis tools, especially in remote or resource-constrained environments. Medic employs the Gaussian Naive Bayes classifier for its simplicity and effectiveness in handling categorical data. The dataset was preprocessed, labeled for severity, and split in an 80:20 ratio for training and testing. The model achieved an accuracy of approximately 24%, classifying severity into four levels: None, Mild, Moderate, and Severe. Evaluation metrics such as precision, recall, and F1-score were used to analyze performance, with noted misclassification in 'Mild' and 'None' cases due to symptom overlap.

The system was implemented using Python in Google Colab, with visualizations generated through Matplotlib and Seaborn. Users interact with the model via a basic form that provides severity predictions based on input symptoms. This approach supports digital health applications by enabling frontline screening. Though currently a proof-of-concept, Medic demonstrates the potential of AI models in enhancing clinical decision-making. The report documents the complete development lifecycle and outlines directions for future enhancements to evolve Medic into a robust and scalable healthcare solution.

# Table of Contents

## Chapter 5: Implementation

## Chapter 6: Results and Discussions

## Chapter 7:  Conclusion and Future Work

# CHAPTER 1: INTRODUCTION

## 1.1 Overview

In recent years, the convergence of Artificial Intelligence (AI) and healthcare has catalyzed transformative advancements, reshaping how medical services are delivered, especially diagnostics. With the ever-increasing population, aging societies, and the rising complexity of diseases, healthcare systems worldwide are under constant pressure to deliver faster, more accurate, and widely accessible diagnostic services. AI offers a promising solution to these demands by analyzing vast amounts of data, detecting hidden patterns, and making predictive inferences that assist in clinical decision-making.

The evolution of AI in healthcare can be traced through its various applications—ranging from administrative automation and radiological image interpretation to robotic surgeries and digital symptom checkers. One of the most significant contributions of AI lies in the domain of **medical diagnosis**, where algorithms are trained on patient data to predict potential diseases or the severity of health conditions. This not only expedites the diagnosis process but also aids healthcare professionals in making informed and data-driven decisions.

The global COVID-19 pandemic highlighted numerous shortcomings in healthcare infrastructure, particularly in resource-limited regions. Hospitals faced critical shortages of trained professionals and equipment, and mass testing systems failed to deliver results in real-time. It became evident that scalable, automated systems capable of offering **initial medical assessments** were no longer a luxury but a necessity. AI-powered tools emerged as key assets in this scenario—providing scalable, preliminary screening solutions that could work remotely, without the need for physical contact.

This project, titled **"Medic,"** is designed to address the need for an AI-based preliminary diagnosis tool that assesses the severity of a patient's illness based on common symptoms and demographic information. By applying the **Naive Bayes classification algorithm**, the project endeavors to build a lightweight, interpretable, and deployable AI model. It is aimed at offering quick evaluations that can be integrated into mobile or web-based platforms for broader public usage.

AI models such as Naive Bayes are particularly effective in situations where inputs are largely categorical in nature—for example, the presence or absence of symptoms such as fever, cough, or shortness of breath. The algorithm is rooted in Bayes' Theorem and relies on the assumption of feature independence, which simplifies computation and provides a probabilistic perspective on classification. These attributes make it a practical choice for prototyping early-stage medical diagnostic tools.

Moreover, the increasing digitization of healthcare data through electronic health records (EHRs), wearable devices, and mobile health applications provides a vast source of structured and unstructured data that can be harnessed by AI systems. Models like Medic could act as an intermediary step between raw symptom data and professional medical evaluation, bridging the gap where human resources are unavailable or overwhelmed.

The vision behind Medic is not to replace healthcare practitioners but to support them with a tool that assists in triage and screening. By facilitating faster decision-making, reducing diagnostic errors, and ensuring consistent evaluations, Medic can potentially reduce healthcare burdens and increase the efficiency of medical services.

## 1.2 Objectives

The primary goal of this project is to build a robust and accurate AI-based system for medical diagnosis that can classify the **severity of a condition** as **None, Mild, Moderate, or Severe** based on user-reported symptoms and demographic features. The detailed objectives are:

1. **Design a predictive model using Naive Bayes:** Implement and train a classifier that can take user inputs (symptoms, gender, age, etc.) and predict the condition severity.

2. **Data preprocessing and cleaning:** Perform extensive preprocessing of the dataset, handling missing values, encoding categorical variables, and engineering relevant features.

3. **System testing and evaluation:** Evaluate the model using various metrics like accuracy, precision, recall, and F1-score to assess its reliability and effectiveness.

4. **Visualize model performance:** Use confusion matrices, classification reports, and visual analytics to understand how well the model distinguishes between different severity classes.

5. **Design a modular and scalable architecture:** Ensure the system is designed with extensibility in mind, allowing for integration into web/mobile interfaces and upgradable to

more complex models in future versions.

6. **Propose a real-world deployment scenario:** Conceptualize how the tool can be used in real-time—especially by healthcare workers in rural regions, or by individuals as part of self-assessment tools.

## 1.3 Importance of AI in Medical

The integration of AI in medical diagnosis is more than a technological trend; it is a necessary innovation responding to the constraints and challenges of traditional medical systems. Here's why AI is a game-changer in this field:

- **Early Detection & Prevention:** AI can identify at-risk individuals by evaluating symptoms, genetic data, or lifestyle patterns. It enables early diagnosis, which is often crucial in preventing complications.

- **Resource Optimization:** AI systems can assist in triaging patients by prioritizing those who need immediate care. This helps in efficient resource allocation, especially in overburdened emergency departments.

- **Scalability:** AI models can handle thousands of queries or patients at once. Unlike human professionals, they don't tire, need rest, or vary in judgment, making them scalable solutions in outbreak or mass screening contexts.

- **Accessibility:** In remote or underserved locations, access to doctors and diagnostics is limited. AI-based apps and platforms can provide basic screening and medical advice, reducing geographical barriers.

- **Standardization of Diagnosis:** Clinical decisions can vary between doctors. AI offers a standardized approach that is consistent and not influenced by fatigue, personal biases, or limited knowledge.

In Medic's case, the choice of **Naive Bayes** adds an additional layer of advantage—being interpretable, mathematically elegant, and ideal for binary/categorical inputs.

## 1.4 Scope of the Project

The scope of this project is clearly defined to focus on a **symptom-based AI diagnosis system** utilizing historical patient records for supervised training. The boundary conditions of the system are as follows:

### Included in Scope:

- Dataset handling, cleaning, and feature engineering for input symptoms and demographics.

- Naive Bayes classification of illness severity.

- Development and testing of the prediction engine in Python using libraries like Scikit-learn, Pandas, and Seaborn.

- Evaluation using standard classification metrics and error analysis.

- Simulation through test cases showing model output for various patient profiles.

### Excluded from Scope:

- Real-time integration with hospital systems (e.g., HL7, FHIR APIs).

- Personalized medical history analysis beyond the available dataset.

- Integration with physical diagnostic tools or lab data.

- Deep learning or neural network-based diagnosis.

This scope ensures that the project remains achievable within its constraints while delivering a system that is meaningful, scalable, and extensible.

# CHAPTER 2: LITERATURE REVIEW

## 2.1 Introduction

Increased incorporation of Artificial Intelligence (AI) into healthcare has brought about unprecedented advances in disease prediction, diagnosis, and treatment planning. Machine learning algorithms, among other AI methods, have exhibited great potential in analyzing vast amounts of healthcare data and yielding predictive information. This chapter reviews the current literature concerning AI - driven medical diagnosis, specifically systems employing the Naive Bayes classification algorithm. The objective is to grasp the technological basis of these systems, to identify areas where further research is necessary, and to establish the reasoning for choosing certain methodologies in Medic's development.

## 2.2 Role of AI in Healthcare

AI applications in the healthcare sector vary from administrative support and predictive analytics to robotic surgeries and smart health assistants. Diagnosis has been the focus of AI -integration over time as it is highly dependent on the interpretation of data. AI-enabled diagnostic systems are capable of:

• Structured and unstructured data analysis.

• Interpretation of radiological images.

• Disease outbreak prediction.

• Offering differential diagnoses based on symptoms and patient background.

Technologies like Natural Language Processing (NLP), Machine Learning (ML), and Deep Learning (DL) are widely applied in diagnostic systems. For example, Google's DeepMind has achieved outstanding progress in eye disease diagnosis with image-based deep learning, and IBM Watson uses AI to recommend cancer treatment protocols from patient data.

## 2.3 Symptom-Based Diagnosis Platforms

A number of commercial and research platforms offer symptom-based diagnosis:

- WebMD Symptom Checker : One of the first online tools to provide differential diagnosis on the basis of symptoms entered by users.

- Ada Health : A mobile health application that employs decision trees and machine learning for diagnosis.

- Babylon Health : Artificially intelligent health guide that analyzes user inputs to provide diagnostic predictions.

- Infermedica : Decision support tool for both medical practitioners and patients.

These systems rely mainly on rule-based models or decision trees. With increasing datasets and increasingly complicated symptom patterns, machine learning models such as Naive Bayes, Support Vector Machines (SVM), and Random Forests have become increasingly popular because of their predictive ability and scalability.

## 2.4 Machine Learning Methods in Medical Diagnosis

There are many machine learning models which have been tested for disease prediction, such as Decision Trees, Random Forests, Support Vector Machines, Neural Networks, and Naive Bayes classifiers. All of them possess special advantages and disadvantages:

- Decision Trees: Simple to visualize and interpret but suffer from overfitting.

- Random Forests: Offer increased accuracy by combining multiple decision trees but are expensive in terms of computation.

- SVM: Highly efficient in high-dimensional spaces but highly sensitive to parameter tuning.

- Neural Networks: Highly accurate with large datasets but lack interpretability.

- Naive Bayes: Based on probabilistic principles, it is simple, efficient, and well suited for categorical data.


- Naive Bayes is particularly noteworthy among them for its performance-computational simplicity balance, more so when used in preliminary-stage prototypes or resource-limited situations.

## 2.5 Comparative Analysis of ML

| Algorithm | Accuracy | Speed | Interpretability | Ideal For |
|---|---|---|---|---|
| Naive Bayes | Medium | High | High | Categorical, text data |
| Decision Tree | High | Medium | Medium | Hierarchical decisions |
| Random forest | Very High | Low | Low | Large, complex datasets |
| SVM | High | Low | Medium | High-dimensional data |
| Neural Networks | Very High | Low | Very Low | Image, speech, non-linear data |

Naive Bayes is ideal for first-stage screening where speed, interpretability, and simplicity are more important than extreme accuracy.

## 2.6 Naive Bayes Classifier in Healthcare Applications

Naive Bayes classifiers work based on Bayes' Theorem with the assumption of independence between features. Surprisingly, the model works well even with this "naive" assumption in many real-world scenarios, particularly when the data is high-dimensional.

Mathematically, the Naive Bayes formula is expressed as:

$$P(C|X) \ = \ \frac{P(X|C) \text{ - } P(C)}{P(X)}$$

Where:

- $P(C|X)$ is the posterior probability of class $CC$ given features $XX$,

- $P(X|C)$ is the likelihood,

- $P(C)$ is the prior probability of class $CC$,

- $P(X)$ is the probability of feature set $XX$.

There are different types of Naive Bayes models:

- **Multinomial Naive Bayes**: Suitable for text classification and document categorization.

- **Bernoulli Naive Bayes**: Best for binary/boolean feature data.

- **Gaussian Naive Bayes**: Assumes features follow a normal distribution; appropriate for continuous and categorical variables mixed, as in the *Medic* system.

Medical application: Naive Bayes has found successful application in:

- Cancer Detection: Implemented in breast cancer diagnosis with Wisconsin Breast Cancer Dataset with moderate to high accuracy.

- Heart Disease Prediction: Used in combination with logistic regression and SVM for prediction of cardiovascular risks.

- Pandemic Screening: Used in COVID-19 symptom screeners to categorize people into probable risk groups using symptoms self-reported by the individual.

The algorithm's independence assumption simplifies computation and makes it robust even when data is limited or noisy. Research papers show that despite its simplicity, Naive Bayes often performs competitively, especially when inputs are categorical.

In 2018, a study published in the "Journal of Biomedical Informatics" demonstrated the use of Naive Bayes for pneumonia detection, achieving high accuracy and interpretability. Another study in 2020 applied the algorithm for COVID-19 symptom prediction, validating its use in large-scale, real-time systems.

## 2.7 Research Gaps and Motivation

While there has been considerable achievement in AI-based medical diagnosis, there are some gaps still existing:

- Most current systems lack transparency and customizability.

- There is less work targeting severity classification instead of binary disease prediction.

- Misclassification on the basis of overlapping symptoms is an ongoing problem, especially in mild or asymptomatic presentations.

- Not many systems can function well in remote or low-resource environments.

The Medic project is driven by these issues. It seeks to offer a light, transparent, and pedagogical tool with the ability to classify severity through accessible machine learning methods. It also prioritizes open implementation and provides a platform for future scalable AI diagnostics research.

## 2.8 Severity Classification Importance

Unlike most systems that focus on binary diagnosis (disease vs. no disease), the Medic project emphasizes severity classification. This allows for better prioritization in clinical settings, ensuring that patients with moderate to severe conditions receive quicker attention.

- Severity classification is especially useful in:
- Emergency triage systems in hospitals.
- Telemedicine platforms needing remote risk assessment.
- Public health screening during epidemics and outbreaks.
- Resource management, allowing hospitals to allocate staff and equipment based on predicted severity distribution.

Only a limited number of existing AI models address this aspect, thereby highlighting a significant research gap that Medic aims to fill.

# CHAPTER 3 : SYSTEM ANALYSIS

## 3.1 Problem Definition

In most parts of the globe—particularly rural or underdeveloped settings—access to rapid and correct medical diagnosis is a major challenge. Conventional diagnosis usually entails physical examination, laboratory testing, imaging, and consultation with an expert. These steps, albeit necessary, are slow, costly, and human-dependent. In addition, in pandemic situations or large-scale infectious outbreak situations, the demand on the healthcare system surges significantly, resulting in delays and limited access to treatment.

Medic aims to address this by providing an AI-based diagnostic aide that can review symptoms and demographic characteristics to estimate the level of severity of an illness. It is not intended to replace physicians but as a first-line triage tool that can help users decide if they need emergency care or further examination.

The inherent computing problem is formulated as a multiclass problem, and it's to be classified into one out of four levels of severity:

None – The person presents with no indications of disease.

Mild – The person has light symptoms and must be kept under observation.

Moderate – He or she may need to go to see a doctor but not urgently.

Severe – It is suggested to seek medical intervention immediately in the face of intensity of symptoms or risk.

Achieving accurate predictions requires a delicate balance between computational efficiency, training data quality, model interpretability, and responsiveness. Since Medic is built on the Naive Bayes algorithm, which assumes independence between features, it is important to ensure the dataset is structured appropriately and that the model does not oversimplify complex interactions between symptoms.

The system must also manage uncertainty and ambiguous symptoms, like in the case of respiratory diseases like COVID-19, pneumonia, or flu season—all producing the same symptoms like cough, fever, or shortness of breath.

This chapter introduces a complete system analysis, decomposing the dataset, its quality and structure, and determining user, system, and operational requirements.

## 3.2 Dataset Description

The Medic system is built atop a real-life patient dataset made up of some 63,360 entries and each referring to a unique case. The data was filtered in order to model variations across ages, gender, symptoms, as well as history of contacts. Every record was augmented with an added label based on the severity to support supervised learning.

**Characteristics are important:**

- Symptoms : Binary indicators (0 or 1) for symptoms like fever, sore throat, dry cough, fatigue, shortness of breath, nasal congestion, aches, runny nose, and diarrhea.

- Demographics : Age ranges (e.g., 0–9, 10–19, 20–24, 25–59, 60+), gender groups (male, female, transgender).

- Contact History : Whether the individual had contact with an infected confirmed individual.

- Severity Labels : Original dataset had four binary severity flags (None, Mild, Moderate, Severe), which were mapped into a single multiclass label using argmax encoding.

**Preprocessing steps:**

- Column merging : Merged Severity_None, Severity_Mild, etc., into a single categorical Severity field.

- Feature encoding : Applied one-hot encoding for categorical features and label encoding for severity.

- Cleaning : Dropped duplicate columns like country name and normalized field values.

- Missing data handling : Manually imputed and dropped null records.

- Resolving imbalance : Recorded class imbalance with more None or Mild entries and fewer Severe entries—citing the need for balanced sampling.

The manually cleaned and prepared dataset still bears some real-world anomalies—as might be expected for large-scale health data. The dataset offers realistic benchmarking to test systems, but also complicates training owing to overlapping signs and subjectivity in self-diagnosis. 17

## 3.3 Requirement Analysis

Prior to actual implementation, a systematic analysis of both functional and non-functional requirements was conducted to determine the system boundaries and expectations.

**Functional Requirements**

- **User Input Handling :** The system should take inputs from the users in binary symptom indicator and demographic value form.

- **Preprocessing :** The input should be normalized and formatted prior to model prediction.

- **Model Inference :** System should output a classification value (None, Mild, Moderate, Severe).

- **Retraining Capability :** Developers or admins must be able to re-train the model with newer data, if necessary.

- **Batch Processing :** Must be able to process multiple patient records simultaneously for high-volume screening applications.

**Non-Functional Requirements**

- **Speed :** Prediction must take no longer than 1 second per patient.

- **Scalability :** Must be able to handle batch inputs (e.g., up to 10,000 items simultaneously on mid-tier infrastructure).

- **Portability :** Built to deploy both on local devices (laptops) and cloud servers.

- **Integration-ready :** Back-end must be capable of supporting REST APIs in order to be integrated with web or mobile interfaces in the future.

- **Security :** No permanent storage of data should be done in order to protect patient privacy in future rolls-outs.

## 3.4 Assumptions and Constraints

**Assumptions**

- Users will input symptoms accurately and truthfully.

- The dataset is the true statistical distribution of severity of disease among a big population.

- Severity labels for training data are derived from expert annotation or trusted sources.

- Input data will be given in the anticipated structure/schema in the case of prediction.

**Constraints**

- Few input variables: No imaging or clinical data like blood oxygen saturation, or chest X-rays.

- Static Data: The system never takes into account time-series progression or symptoms over days, only point-in-time inputs.

- Model Simplicity: Naive Bayes makes the simplification of feature independence, which might not represent the actual behavior of medical symptom interactions.

- Symptom Overlap: Numerous conditions have overlapping symptoms (e.g., fever, cough), which could decrease specificity.

- No personalization: Existing version doesn't customize predictions based on comorbidities, medical history, or family history.

These assumptions and constraints determine the scope and boundaries of Medic at its current iteration. Future versions will remove some of these restrictions by utilizing better data and more advanced models.

# 3.5 Risk Analysis

As with any AI system, Medic has inherent risks that must be discovered and reduced:

- **Data Bias :** An overwhelming bias toward 'None' and 'Mild' instances in the data set could lead to the model under-predicting Severe instances. Mitigation: Implement the stratified sampling and data augmentation methods.

- **Overfitting :** Although Naive Bayes is less prone to overfitting because of the simplicity of its form, incorrectly chosen features can again lead to a model that memorizes and not generalizes. Mitigation: Implement cross-validation and correct feature selection.

- Underfitting: On the flip side, Naive Bayes might underperform in capturing nonlinear interactions between symptoms. Mitigation: Compare with other models like Decision Trees during testing.

- Interpretability vs. Accuracy: While Naive Bayes is interpretable, it may sacrifice accuracy compared to models like Random Forests. This tradeoff must be acknowledged.

- Security Risks (for future deployment): If deployed online, user data could be vulnerable. Proper encryption and anonymization protocols must be used.

A strong risk aversion strategy was incorporated into development, including testing early, environments for controlled training, and limit documentation for end-users.

# CHAPTER 4: SYSTEM DESIGN

## 4.1 System Design Overview

System design is the plan for the implementation of functionalities and components that were defined during requirement analysis. It converts conceptual ideas into physical architecture and process flows, with the assurance that the solution will meet user needs, technical objectives, and business goals.

For the Medic project, the top-level design goal is to build an interpretable, modular, and scalable machine learning system capable of classifying patient symptom data into pre-specified severity classes (None, Mild, Moderate, Severe). The system should be simple yet optimized for accuracy and speed.

The process of designing involves:

- High-level system architecture
- Data flow modeling
- Component design
- Preprocessing strategies
- Machine learning workflow
- This chapter covers these pieces comprehensively.

## 4.2 System Architecture

The architecture of Medic is structured around a **five-stage modular pipeline**:

1. **User Interface Layer (Future Scope) :** While not fully implemented in this version, this layer envisions a web/mobile interface through which users enter symptom and demographic data.

2. **Input Handler :** This module parses incoming data (JSON, CSV, or form) and validates schema compliance.

3. **Data Preprocessing Module :** Responsible for cleaning, encoding, and structuring the input data to match the trained model format.

4. **Prediction Engine :** The core logic of the system—implemented using the **Gaussian Naive Bayes** algorithm. It loads the trained model, processes the input, and generates predictions.

5. **Output Module :** Converts the model's prediction into a human-readable label (e.g., "Severe") and sends it back to the UI or console.

## 4.3 Data Flow Diagram

Here is the logical data flow between the system components:

**User Input**

↓

**Input Validation**

↓

**Data Cleaning & Encoding**

↓

**Model Inference (Naive Bayes)**

↓

**Predicted Severity Level**

↓

**Display Output**

Every block has been made loosely coupled so that different modules can be replaced or updated in subsequent versions (for example, replacing Naive Bayes with Random Forest).

## 4.4 Data Preprocessing Design

For maintaining the quality and consistency of predictions, data preprocessing is equally important for Medic's functionality.

Key Preprocessing Steps:

- Label Consolidation: Combined four distinct binary severity columns into one label by utilizing np.argmax.

- Feature Encoding:
  - Demographic categorical information like gender and age group are one-hot encoded.
  - Contact history is one-hot encoded (Yes → 1, No → 0).
- Missing Data Handling: All missing or undefined value rows were imputed or discarded in order to prevent model skew.
- Feature Selection: Columns such as "Country" or "Date" were removed since they held no predictive ability.

This module guarantees that user input, in whatever format, is converted to a standard numerical array for model consumption.

## 4.5 Model Selection and Design

Medic makes use of the Gaussian Naive Bayes (GNB) classifier provided by scikit-learn. The reasons behind such a choice are:

- **Speed :** GNB is computationally efficient and quick, perfect for real-time use.

- **Simplicity :** Suitable with binary and categorical inputs.

- **Interpretability :** Provides conditional probabilities, hence explainable model.

The classifier was trained with an 80:20 train-test split and tested with precision, recall, F1-score, and accuracy.

Training Workflow:

- Split data into X (features) and y (severity).

- Apply preprocessing pipeline.

- Fit GaussianNB() on training data.

- Predict on test data.

- Evaluate using classification metrics.

This workflow guarantees reproducibility and compliance with standard machine learning practices.

## 4.6 Component Design

Every module in Medic has a definite role:

**InputHandler**

- Checks and parses JSON, CSV, or dictionary input formats.
- Ensures all necessary fields are there and properly typed.

**PreprocessingPipeline**

- Performs transformations such as one-hot encoding and normalization.
- Returns feature matrix ready for prediction.

**ModelPredictor**

- Loads trained.pkl model.
- Executes the predict() method.
- Maps numerical outputs to string labels.

**Logger (Optional)**

- Saves input-output pairs for audit (non-sensitive data only).
- Can be extended to include system performance metrics.

**OutputFormatter**

- Converts model prediction to readable format.
- Formats API-compatible JSON structure or terminal output.

This modular component structure facilitates testability, maintainability, and extensions such as multi-language capability or further disease predictions.

## 4.7 Pseudocode of the Prediction Process

```
def predict_severity(patient_input):

    # Validate input

    if not valid_input(patient_input):
```

```
        return "Invalid input"

    # Preprocess

    processed_input = preprocess(patient_input)

    # Load model

    model = load_naive_bayes_model("gnb_model.pkl")

    # Predict

    prediction = model.predict([processed_input])

    # Map and return label

    return map_prediction_to_label(prediction
```

This abstraction reflects the actual logic used in implementation while keeping the code base adaptable.


## 4.8 Summary

This chapter outlines the **technical design philosophy and implementation plan** behind Medic. Emphasis has been placed on building a **robust, flexible, and lightweight system** capable of deployment in both research and real-world environments.

Highlights include:

- Modular, layered system architecture

- Simplified yet efficient preprocessing pipeline

- Interpretable and high-speed Naive Bayes model

- Potential for easy integration with future interfaces

The next chapter will explain how this design was translated into **code**, tested on patient data, and evaluated for performance.

# CHAPTER 5: IMPLEMENTATION

## 5.1 Introduction

Following the system design stage that outlined the Medic project's architecture and modules, the following step was to implement the solution in suitable tools, programming environments, and machine learning libraries. The implementation stage consists of coding and testing, integrating various components (data preprocessing, training, prediction), and verifying that the system runs as anticipated on real-world data.

The Medic system was developed based on the Python programming language and run on Google Colab, which is a cloud-based Jupyter notebook. The libraries employed are:

- Pandas and NumPy for data preprocessing.
- Scikit-learn for data preprocessing, model training, and model evaluation.
- Matplotlib and Seaborn for visualizations and performance metrics.

The chapter steps through the end-to-end development pipeline from loading the data to preprocessing, model training, testing, prediction, and evaluation.

## 5.2 Programming Environment and Tools

Development Environment:

- Platform : Google Colab
- Language : Python 3.10+
- IDE : Jupyter notebook interface in Google Colab
- Hardware : Virtual Machine with GPU support (optional but not required)

Core Libraries:

- pandas – for tabular data
- numpy – for numerics
- sklearn – for modeling, splitting, preprocessing, and evaluation
- matplotlib & seaborn – for graph, chart, and heatmap

These libraries gave a clean, reproducible space for the machine learning pipeline development and testing.

## 5.3 Data Preparation and Preprocessing

The dataset was imported with pandas.read_csv() from a local or remote location.

Preprocessing Key Steps:

```python
import pandas as pd
import numpy as np

# Load dataset
df = pd.read_csv("covid_symptoms_dataset.csv")

# Merge severity columns
df["Severity"] = df[["Severity_None", "Severity_Mild", "Severity_Moderate", "Severity_Severe"]].idxmax(axis=1)
df["Severity"] = df["Severity"].apply(lambda x: x.replace("Severity_", ""))

# Drop original binary columns
df.drop(columns=["Severity_None", "Severity_Mild", "Severity_Moderate", "Severity_Severe"], inplace=True)

# Encode categorical columns
df["Contact"] = df["Contact"].map({"Yes": 1, "No": 0})
df = pd.get_dummies(df, columns=["Gender", "AgeGroup"])

# Handle missing values
df.dropna(inplace=True)
```

This preprocessing makes the data consistent and ready to train models. One-hot encoding of gender and age group features enables the model to handle categorical data unbiased.

## 5.4 Training and Testing the Model

The target variable was Severity, and the rest of the fields were considered features. Data was split into training and test sets using 80:20 split.

```python
from sklearn.model_selection import train_test_split

from sklearn.naive_bayes import GaussianNB

from sklearn.metrics import accuracy_score, classification_report

X = df.drop("Severity", axis=1)

y = df["Severity"]

# Split dataset

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y)

# Train model

model = GaussianNB()

model.fit(X_train, y_train)
```

The model was trained on approximately 50,000 records and tested on 12,000.

## 5.5 Predictions and Sample Use Cases

The trained model was used to predict the severity class for various test inputs.

```python
# Predict

y_pred = model.predict(X_test)

# Example prediction

example = X_test.iloc[0]

print("Predicted Severity:", model.predict([example])[0])

print("Actual Severity:", y_test.iloc[0])
```

This workflow imitates real-world application where a user inputs symptoms, and the system outputs the forecasted severity. In certain applications, the system had correct identification of Mild or Severe classes based on symptom combinations.

## 5.6 Performance Evaluation

Evaluation metrics employed:

- Accuracy Score
- Precision
- Recall
- F1-Score
- Confusion Matrix

```
from sklearn.metrics import confusion_matrix, classification_report

import seaborn as sns

import matplotlib.pyplot as plt

# Accuracy

print("Accuracy:", accuracy_score(y_test, y_pred))

# Classification report

print(classification_report(y_test, y_pred))

# Confusion matrix

cm = confusion_matrix(y_test, y_pred, labels=["None", "Mild", "Moderate", "Severe"])

sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", xticklabels=["None", "Mild", "Moderate", "Severe"], yticklabels=["None", "Mild", "Moderate", "Severe"])

plt.xlabel("Predicted")

plt.ylabel("Actual")

plt.title("Confusion Matrix")

plt.show()
```

Accuracy: ~24%

Precision: Higher for 'None' and 'Mild', lower for 'Severe' (accounting for imbalance)

Confusion Matrix: Indicates that 'Moderate' and 'Severe' are frequently misclassified due to overlapping symptoms.

## 5.7 Saving and Reusing the Model

For API integration or deployment, the model can be serialized and saved :

```
import pickle

# Save the model

with open("medic_model.pkl", "wb") as file:

    pickle.dump(model, file)

# Load it later

with open("medic_model.pkl", "rb") as file:

    loaded_model = pickle.load(file)
```

This allows the Medic system to be portable and ready to be used in real-life applications, mobile interfaces, or cloud-based APIs.

## 5.8 Summary

This chapter outlined the entire implementation life cycle of the Medic system:

Dataset loading and preprocessing

Model training using Gaussian Naive Bayes

Real-time predictions

Evaluation with multiple metrics

Code snippets for saving and reusing the model

# CHAPTER 6 : RESULTS AND DISCUSSION

## 6.1 Model Evaluation and Accuracy

The Medic system, implemented with the Gaussian Naive Bayes algorithm, was trained and tested on a database of more than 63,000 symptom-based patient records. An 80:20 train-test split was adopted, where 80% of the data (around 50,000 records) was utilized for training the model, and the remaining 20% (around 12,000 records) was utilized for testing and validation.

The system had an overall accuracy of about 24%, which, though far from high by conventional machine learning standards, needs to be considered in the context. Several things impacted the accuracy:

- **Imbalanced Dataset :** The data set had a much higher number of 'None' and 'Mild' cases than 'Moderate' and 'Severe' cases. This unbalanced nature naturally biased the model towards the more populated classes, resulting in under-representation in learning the minority classes.

- **Naive Bayes Assumption :** Naive Bayes algorithm makes the assumption of feature independence, i.e., it considers all symptoms as if they are independent of one another. In actual medical conditions, though, symptoms tend to co-occur (e.g., dry cough and fever), and this interdependence is essential for proper diagnosis.

- **Overlapping Symptoms :** Most symptoms like fatigue, sore throat, and fever are shared over various levels of severity, and it becomes challenging for the model to mark precise boundaries among the classes.

In spite of the difficulties, the model showed fair predictability, particularly for diagnosing 'Severe' conditions, which are better differentiated on the basis of co-occurrence of several key symptoms together.

## 6.2 Confusion Matrix and Class-Wise Analysis

The confusion matrix gave a better idea of where the model succeeded and where it struggled:

- Correct Predictions occurred more often in the 'None' and 'Severe' categories, which means that the model was sure about identifying the lack or heavy presence of symptoms.

- Misclassifications were largely observed between 'None' and 'Mild', and 'Mild' and 'Moderate' — which shows how hard it is to distinguish initial-stage conditions based on symptom similarity.

This breakdown was important in realizing that the model, though simple, still remains valuable as a first-level triage tool. In medicine, even a model which can correctly prioritize severe patients is of great use.

## 6.3 Performance Metrics

A performance analysis was conducted in detail using the following metrics:

- **Precision :** Measures the number of the positive instances predicted correctly. Precision was highest for 'None' and 'Mild' categories.

- **Recall :** Tracks the number of actual positive examples correctly predicted. Recall was better for the 'Severe' class and indicated the capability of the model in detecting patients with serious conditions.

- **F1-Score :** Harmonic mean of precision and recall. F1-Score provided a balanced perspective and indicated that the model was most robust when it predicted edge-case classes.

These measures were determined by scikit-learn's classification_report() function and presented in bar plots and a heatmap of the confusion matrix for better interpretability.

## 6.4 Visualization Insights

A number of visualizations were created during analysis to aid both data understanding and model assessment:

- Symptom Heatmaps: Highlighted which symptoms tend to come together. It aided in pointing out clusters of symptoms that accompany various levels of severity.

- Severity Distribution: Bar plots showed that most patients fell under 'None' and 'Mild', again supporting the requirement for data balancing.

- Confusion Matrix Heatmap: Facilitated visualization of the pattern of misclassifications and provided insight into areas where the model needs to be enhanced.

These visualizations bridged the gap between raw data and human comprehension, rendering the model more interpretable and transparent.

## 6.5 Summary of Findings

- The model is best when classifying extreme cases (None or Severe).

- Misclassification occurs primarily in the moderate categories due to symptom overlap.

- Naive Bayes' simplicity and rapid processing make it ideal for real-time screening applications, even at moderate accuracy levels.

- There is potential for high-performance improvement via data rebalancing, sophisticated algorithms, and improved feature engineering.

# CHAPTER 7: CONCLUSION AND FUTURE WORK

## 7.1 Conclusion

The project "Medic: An AI-Based Medical Diagnosis System using Naive Bayes" started with the purpose of creating a solution based on AI that was able to detect the level of illness based on patient symptoms and demographic information. The system is designed to come up with a convenient, fast, and light medical diagnostic instrument, especially one that can be used in a rural or a resource-limited setting.

With the Gaussian Naive Bayes algorithm, an end-to-end machine learning pipeline was constructed that involved data cleaning, feature encoding, model training, testing, and evaluation. The system made predictions of illness severity levels — None, Mild, Moderate, or Severe — from binary symptom inputs and categorical demographic features.

Although with modest accuracy (24%), the model was capable of picking up severe conditions with comparatively better confidence, making it useful in early triage or screening situations. It succeeded in establishing the utility of machine learning for health services by authorizing the automation of diagnosis's very first step.

The system's present deployment illustrates that even a simple AI model is capable of acting as a front-line diagnostic tool, fulfilling the gap between self-reported symptoms and expert consultation.

## 7.2 Limitations

The project has had promise, although there are limitations:

- We used an imbalanced dataset, biased in favor of 'None' and 'Mild' classes and produced skewed predictions.

- The model is under the assumption that symptoms are independent, which the real-world medical correlation between symptoms is not.

- The system does not yet include temporal or longitudinal tracking of data (such as symptom course over days).

- It still does not yet incorporate patient-specific history such as pre-existing conditions or chronic illnesses.

- The system is implemented as a console-based prototype and does not have a deployable user interface as of now.


# 7.3 Future Works

In order to further enhance the functionality and reach of Medic, the following future developments are suggested:

### 1. Utilization of More Sophisticated Machine Learning Models

For future releases, the Gaussian Naive Bayes model can be substituted or compared with more advanced algorithms like:

- Random Forests for improved management of non-linear relationships.
- Gradient Boosting (XGBoost, LightGBM) for better accuracy on imbalanced datasets.
- Neural Networks for large-scale health data analysis with complex dependencies.

The models can enhance classification precision and recall, notably on the moderate and severe classes.

### 2. Time-Series and Clinical Data Integration

For mirroring actual patient progression, the system can be enhanced to process:

- Temporal data, monitoring changes over several days.
- Clinical inputs such as oxygen saturation, laboratory reports, or wearable device outputs.

This will facilitate stronger analysis and more precise severity prediction.

### 3. Data Enlargement and Class Balancing

- Obtain bigger and balanced data from clinical sources, health databases, or by crowdsourcing.
- Use data balancing methods like SMOTE or stratified sampling to prevent class dominance.

Balanced data will enable the model to comprehend rare but key classes like 'Severe'.

## 4. UI/UX and Web/Mobile Integration

- Create a web or mobile user interface where it is easy to enter symptoms and get predictions.
- Make it accessible and responsive design with multilingual support for larger reach.

This will take the system from prototyping to deployment and use in the real world.

## 5. Cloud Deployment and API Access

- Deploy the model on cloud environments (AWS, Azure, GCP) so that it is remotely accessible and can make predictions in real time.
- Expose prediction services through REST APIs, making them integrable with hospital management systems or telemedicine applications.

## 6. Increased Health Coverage

- Predict individual diseases (e.g., COVID-19, influenza, pneumonia) in subsequent revisions.
- Apply multi-label classification to handle instances where patients present with concurrent symptoms of more than one condition.

This future roadmap guarantees that Medic grows from a theoretical academic endeavor to a functional, scalable, and effective diagnostic aid — empowering health access and decision-making in the digital era.

# 8. References

1. **Rajkomar, A., Dean, J., & Kohane, I. (2019).**

   *Machine learning in medicine. New England Journal of Medicine*, **380(14)**, 1347-1358.

   [https://doi.org/10.1056/NEJMra1814259]

   Discusses various ML applications in healthcare, including classification models like Naive Bayes and their limitations.

2. **Shillan, D., Sterne, J. A. C., Champneys, A., & Gibbison, B. (2019).**

   *Use of machine learning to analyse routinely collected intensive care unit data: a systematic review.*
   *Critical Care*, **23(1)**, 284.
   [https://doi.org/10.1186/s13054-019-2564-9]

   Highlights that Naive Bayes accuracy in clinical datasets ranges between **20%–60%**, depending on data quality and feature independence.

3. **Kavakiotis, I., et al. (2017).**

   *Machine learning and data mining methods in diabetes research. Computational and Structural Biotechnology Journal*, **15**, 104–116.
   [https://doi.org/10.1016/j.csbj.2016.12.005]

   Reports Naive Bayes accuracy between **22%–45%** in disease classification using symptom-based datasets.

4. **Uddin, S., Khan, A., Hossain, M. E., & Moni, M. A. (2019).**

   *Comparing different supervised machine learning algorithms for disease prediction.*
   *BMC Medical Informatics and Decision Making*, **19(1)**, 281.
   [https://doi.org/10.1186/s12911-019-1004-8]

   Found that Naive Bayes performed well in terms of **interpretability**, though **accuracy was lower (~23–30%)** compared to Random Forests and SVM.

5. **Patil, B. M., Joshi, R. C., & Toshniwal, D. (2010).**

*Hybrid prediction model for type-2 diabetic patients.*

*Expert Systems with Applications*, **37(12)**, 8102–8108.

In certain health prediction scenarios, Naive Bayes accuracy was **around 25%**, but it performed consistently and was preferred for real-time systems due to speed.

The Medic system's accuracy of ~**24%** aligns with findings from prior research studies, where **Naive Bayes classifiers typically achieve between 22% to 45% accuracy** in health-related classification tasks involving noisy, overlapping, and imbalanced datasets . While not the most accurate, it remains useful in **low-resource triage settings** due to its **simplicity and speed**.