

Dataset Creation and Problem Formulation

Title: Hospital Length of Stay Prediction

1. Problem Statement:

Hospitals face increasing pressure to manage resources efficiently while delivering quality patient care. One major challenge is accurately predicting the length of hospital stay (LOS) for patients at the time of admission. Incorrect estimation of hospital stay can lead to overcrowding, inefficient bed management, increased operational costs, and delayed treatment for other patients.

The length of stay varies based on multiple factors such as patient age, diagnosis, severity of illness, number of procedures, and existing medical conditions. Currently, many hospitals rely on doctors' experience or historical averages, which may not always reflect individual patient complexity.

This problem is significant because better prediction of hospital stay enables hospitals to optimize bed allocation, reduce waiting time, improve patient flow, and control healthcare costs. It also supports better discharge planning and improves patient satisfaction.

The primary beneficiaries of this solution are hospital administrators, healthcare professionals, and patients. Administrators can manage hospital resources more effectively, doctors can plan treatments better, and patients receive timely care without unnecessary delays.

Thus, predicting hospital length of stay using data-driven approaches is a meaningful and impactful real-world problem suitable for machine learning and data analysis.

2. Proposed Solution Approach

This problem can be addressed using supervised machine learning, specifically regression techniques, because the target variable (length of stay) is a continuous numerical value measured in days.

A regression model such as Multiple Linear Regression, Decision Tree Regressor, or Random Forest Regressor can be trained using historical hospital data. The input features would include patient age, gender, diagnosis, admission type, number of tests, number of procedures, and comorbidities. The model learns the relationship between these features and the length of hospital stay.

The dataset will first be cleaned and preprocessed by handling missing values, encoding categorical variables, and normalizing numerical data if necessary. Then, the dataset will be split into training and testing sets. The model will be trained on the training data and evaluated using performance metrics such as Mean Absolute Error (MAE) and R-squared score.

The expected outcome is a predictive model that can estimate the number of days a patient is likely to stay in the hospital upon admission. This insight allows hospitals to better manage bed occupancy, staff scheduling, and treatment planning.

In addition, explainable machine learning techniques can be applied to understand which factors most strongly influence hospital stay, making the system transparent and trustworthy for medical professionals.

3. Dataset Description:

The dataset used in this project is a self-created dataset representing hospital patient records. It contains information about patients admitted to a hospital and the corresponding length of stay. The dataset consists of 40 records and 9 attributes.

Column Name	Data Type	Description	Relevance
Patient_ID	Integer	Unique identifier for each patient	Helps distinguish individual records
Age	Integer	Age of the patient in years	Older patients often need longer care
Gender	Categorical	Male/Female	Certain conditions vary by gender
Diagnosis	Categorical	Main illness or condition	Strongly affects treatment duration
Admission_Type	Categorical	Emergency or Planned	Emergency cases often stay longer
No_of_Tests	Integer	Number of medical tests	Indicates severity and complexity
No_of_Procedures	Integer	Number of procedures performed	More procedures usually mean longer stay
Comorbidities	Integer	Count of additional diseases	More comorbidities increase risk
Length_of_Stay	Integer	Number of days stayed (Target)	Output variable to predict

The dataset is structured in a tabular format, where each row represents a patient and each column represents a specific feature. Numerical and categorical data types are both present, making the dataset suitable for supervised learning.

This data is necessary and sufficient for the proposed solution because it captures both demographic and clinical factors that directly influence hospital stay duration. Together, these features provide a comprehensive view of patient condition and treatment complexity, enabling accurate prediction.

4. Data Source Documentation

The dataset used in this project is an original, self-created dataset designed for academic and learning purposes. The data was generated manually based on realistic hospital scenarios and domain knowledge about healthcare operations. Values were created to reflect practical relationships, such as higher hospital stay for emergency admissions and patients with more comorbidities.

No real patient data was used, ensuring that there are no privacy or ethical concerns. The dataset simulates real-world hospital data patterns and is suitable for demonstrating data analysis and machine learning techniques in an educational setting.