

## Introduction

Efficient hospital bed management plays a critical role in streamlining healthcare operations, minimizing patient wait times, and optimizing the use of limited medical resources. In response to this operational challenge, this project leverages Natural Language Processing (NLP) and Machine Learning (ML) techniques to analyze clinical notes and predict patients' Length of Stay (LOS) in hospitals.

The project is centered around three core tasks:

1. **Entity Extraction:** Automatically identifying key medical entities from free-text clinical notes, including patient demographics, symptoms, prior conditions, prescribed medications, and medical procedures.
2. **Entity Standardization:** Converting non-standardized medical terms into consistent, structured formats using controlled vocabularies or mapping strategies.
3. **LOS Prediction:** Building a machine learning model to predict the expected number of days a patient will stay in the hospital, based on the clinical information documented during the first day of admission.

## Problem Statement

Hospitals often struggle with accurate forecasting of patient bed occupancy due to diverse patient conditions and inconsistently documented clinical notes. While clinical notes contain a wealth of critical information, their unstructured and variable nature makes it difficult to extract and utilize this data for operational planning.

The primary objectives of this project are:

- To convert unstructured nursing notes into structured, analyzable data.
- To standardize extracted medical entities to ensure consistency across records.
- To develop a robust predictive model that estimates patient LOS using both structured and unstructured data sources.

By achieving these goals, the project aims to support hospital administrators in making data-driven decisions regarding bed allocation and resource management, thereby helping reduce overcrowding and improve patient care outcomes.

## Dataset Overview

The dataset comprises day-wise nursing notes for individual hospital admissions. Each record represents a unique patient encounter and includes both administrative metadata and unstructured clinical notes recorded over multiple days. These notes provide detailed documentation of a patient's condition, treatments, symptoms, and overall progress during their hospital stay.

Some entries in the dataset are marked with "\$\$", which indicates missing nursing notes for that particular day. Handling these missing values appropriately is an important step in the preprocessing pipeline.

## Features

The main features in the dataset include:

- **SUBJECT\_ID**: Unique identifier for each patient.
- **HADM\_ID**: Unique identifier for each hospital admission.
- **ADMITTIME**: Timestamp indicating when the patient was admitted.
- **DISCHTIME**: Timestamp indicating when the patient was discharged.
- **day\_1 to day\_100**: Daily nursing notes starting from the day of admission (day\_1) up to the 100th day (day\_100). Each column contains free-text clinical observations and notes for that day. Missing notes are denoted by "\$\$".

These daily notes form the primary unstructured data source used for entity extraction and LOS prediction. The number of days with notes can vary across patients, reflecting their actual hospital stay durations.

## Methodology

### 1. Named Entity Recognition (Entity Extraction)

The first step in the methodology involves extracting meaningful clinical entities from unstructured nursing notes. This is achieved through Named Entity Recognition (NER) using a domain-specific large language model. The goal is to convert free-text clinical notes into structured, analyzable data formats.

#### 1.1 Data Loading and Preprocessing

- The nursing notes are stored in an Excel file with day-wise documentation (e.g., day\_1, day\_2, ..., day\_100).
- Two sheets were loaded: one containing daily nursing notes, and another containing additional patient details.
- All daily note columns were identified and concatenated into a single text field, combined\_notes, per patient. The function clean\_note\_simple() was used to:
  - Replace masked identifiers (e.g., [\*\*Name\*\*]) with readable tokens.
  - Remove missing note markers (\$\$).
  - Normalize whitespace.

## 1.2 Model Selection and Setup

To perform biomedical entity extraction, the **Bio-Medical-Llama-3-8B** model from Hugging Face was used:

- **Model Loading:**
  - Loaded using transformers with 4-bit quantization (via BitsAndBytesConfig) to optimize memory efficiency and inference speed.
  - Model and tokenizer were loaded with AutoModelForCausalLM and AutoTokenizer, respectively.
  - The model was deployed on GPU using Colab with automatic device mapping.
- **Hugging Face Authentication:**
  - Secure access was granted using a personal access token to load the model from Hugging Face's model hub.

## 1.3 Prompt Engineering

A prompt template was designed to instruct the language model on what to extract:

plaintext

From the clinical note below, extract the following entities and return them in JSON format:

- **PatientDemographics:** Age and Gender.
- **CurrentConditions:** Current diseases or symptoms.
- **PastConditions:** Historical or resolved conditions.
- **Medications:** Prescribed drugs.
- **Procedures:** Surgeries or interventions.

The prompt is injected with the patient's combined\_notes and then tokenized as input for the model.

## 1.4 Entity Extraction Logic

- The extract\_entities() function handles the full pipeline for a single note:
  - Constructs the prompt using build\_prompt().
  - Truncates and tokenizes the input.
  - Generates a response from the model using controlled decoding parameters (max\_new\_tokens, temperature, eos\_token\_id).
  - Decodes the model's response to extract a structured JSON-like output.

## 1.5 Batch Processing and Output

- To process multiple patient records efficiently, the `progress_apply()` function from `tqdm` was used to visualize progress while applying `extract_entities()` across the dataset.
- The results were stored in a new column, `extracted_entities`.
- Finally, the extracted data was saved to an Excel file (`extracted_entities_200_2.xlsx`) for further analysis or downstream tasks.

## 1.6 Hyperparameter Settings

The entity extraction process utilizes a causal language model with carefully chosen decoding hyperparameters to balance response coherence and diversity:

- **max\_new\_tokens=512**: Limits the number of generated tokens to prevent overly long or verbose outputs while ensuring adequate space for all entities.
- **temperature=0.5**: Controls randomness in generation; a moderate value encourages factual and consistent outputs while allowing some flexibility in phrasing.
- **do\_sample=True**: Enables stochastic sampling, making the model outputs more natural rather than deterministic.
- **eos\_token\_id**: Ensures the model stops generating once an end-of-sequence token is reached, improving efficiency and output formatting.

Additionally, the model is loaded using **4-bit quantization** (`load_in_4bit=True`) to significantly reduce memory usage and accelerate inference without sacrificing accuracy — ideal for large-scale processing in a resource-constrained environment like Google Colab.

## 2. Entity Standardization

Following entity extraction, a **standardization** step was performed to map the extracted clinical terms to standardized concepts in the **Unified Medical Language System (UMLS)**. This step ensures interoperability, consistency, and downstream compatibility with structured clinical ontologies. Below is a summary of the methodology and its components:

### Objective:

Map extracted free-text clinical terms (from `CurrentConditions`, `PastConditions`, `Medications`, and `Procedures`) to UMLS Concept Unique Identifiers (CUIs).

## Key Steps:

### 1. UMLS-Compatible Term Extraction

- A prompt-based language model is used to extract **concise, clinically meaningful phrases** from raw free-text.
- The model is the same **Bio-Medical-LLaMA-3-8B** used in the entity extraction phase, leveraging its domain knowledge to return **comma-separated** clinical terms optimized for UMLS lookup.

**prompt** = f"""Extract standardized medical terms suitable for UMLS lookup from the Input. Use short, clinically relevant phrases such as anatomical structures, conditions, or procedures. Return only a comma-separated list of terms (no new line characters \n)

Examples:

Input: C7 Pedicle fracture that extends to facet

Output: pedicle fracture, facet joint

Input: Large disc protrusion with cord compression

Output: disc protrusion', cord compression

Input: ACDF with allograft plate C6-7

Output: ACDF, allograft plate, C6-7

Input: {note}

Output:

""

### 2. UMLS CUI Mapping

- An API-based method connects to the **UMLS REST service**.
- Authentication is handled via a **Ticket Granting Ticket (TGT)** and **Service Tickets (ST)**.
- Each cleaned term is searched against the UMLS database using the search/current endpoint.
- The most relevant **CUI (Concept Unique Identifier)** is retrieved for valid matches.

### Fallback and Handling:

- Terms longer than 20 characters are first cleaned using the LM-based term extractor before UMLS lookup.

- Terms that yield **no result** are skipped.
- A JSON dictionary is constructed, appending CUIs under keys like "CurrentConditions\_CUIs", "Procedures\_CUIs", etc.

### 3. Integration and Output

- Standardization is applied to all the data points.
- Output is saved as a new column: "standardized\_entities" containing a JSON-formatted string with mapped CUIs.
- Final output is exported as an Excel file for further analysis.

#### Output:

```
{"CurrentConditions_CUIs": ["C0877393", "C0238115"], "PastConditions_CUIs": [], "Medications_CUIs": ["C0033487", "C0026549", "C0700776", "C3872070", "C0006699"], "Procedures_CUIs": ["C0877393", "C0014872"]}
```

## LOS Category Prediction – Classification Task Report

This section outlines the **Length of Stay (LOS) category prediction** process, where the goal was to classify patients into "short" ( $LOS \leq 3$  days) or "long" ( $LOS > 3$  days) categories based on structured and semi-structured clinical data. The complete pipeline spans data cleaning, transformation, embedding generation using **Bio\_ClinicalBERT**, and classification using **XGBoost**.

### A. Data Cleaning and Preprocessing

#### 1. Parsing and Structuring Extracted Entities

- JSON-like strings from the extracted\_entities column were parsed using a custom function to remove markdown formatting and convert strings into Python dictionaries.
- Key fields extracted:
  - Age, Gender, CurrentConditions, PastConditions, Medications, Procedures

#### 2. Data Merging

- Patient demographics were merged with other relevant data (e.g., diagnosis, LOS) using HADM\_ID as the join key.

#### 3. Column Cleaning

- Converted list-style string columns (e.g., ['condition1', 'condition2']) into clean comma-separated strings by removing brackets and quotes.
- Defined a custom function to categorize **LOS** into:

- short: 0–3 days
- long: >3 days

#### 4. Age Standardization

- Handled mixed types (strings, ints) and removed inconsistent formats like "Older adult".
- Mapped outliers (age  $\geq 100$ ) to a normalized value (e.g., 65), assuming possible data entry errors.

#### 5. Gender Standardization

- Mapped a variety of noisy gender labels (e.g., "inferred as male", "likely female") into consistent "M", "F", or "Unknown".

### B. Feature Engineering

Combined all textual clinical information into a single input:

#### C. Embedding Generation with Bio\_ClinicalBERT

##### 1. Model Selection

- Used domain-specific pre-trained model: emilyalsentzer/Bio\_ClinicalBERT, designed for clinical text.

##### 2. Embedding Strategy

- For each combined clinical text, extracted the **[CLS] token embedding** from the last hidden layer as a 768-dimensional representation of the entire input.
- This embedding acts as a dense feature vector for downstream classification.

### D. Classification and Evaluation

#### 1. Model:

- Applied **XGBoost**, a robust gradient boosting classifier.

#### 2. Training and Evaluation:

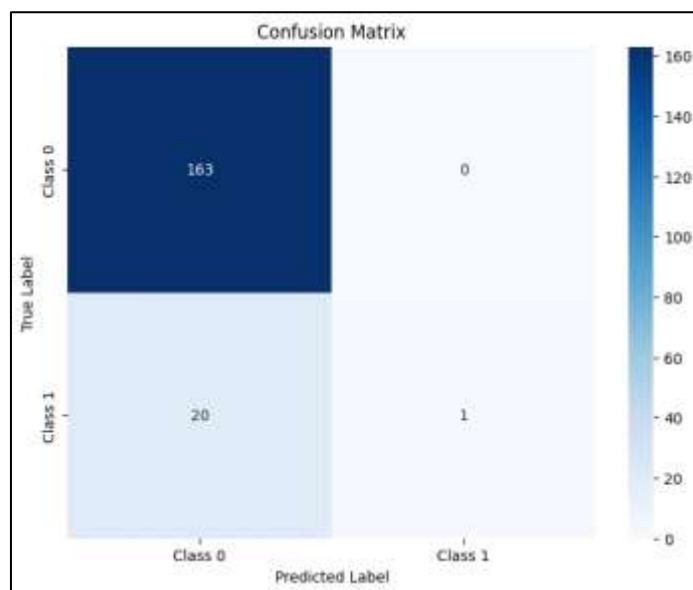
- 80-20 train-test split.
- Labels encoded via LabelEncoder (short  $\rightarrow$  0, long  $\rightarrow$  1).
- Evaluated using:
  - **Precision**
  - **Recall**
  - **F1-Score**

- Accuracy

### 3. Advantages of this setup:

- Bio\_ClinicalBERT provides semantically rich embeddings of complex clinical text.
- XGBoost efficiently handles structured representations and performs well on tabular tasks.

The classification results show **high overall accuracy (89%)**, but a **clear imbalance in performance between the two classes**, especially poor performance on the minority class (short stays). Here's a breakdown:



Result Summary:				
	precision	recall	f1-score	support
False	0.89	1.00	0.94	163
True	1.00	0.05	0.09	21
accuracy			0.89	184
macro avg	0.95	0.52	0.52	184
weighted avg	0.90	0.89	0.85	184



## Conclusion

In this project, our primary objective was to **maximize the precision for predicting long hospital stays**, based on a specific business requirement that prioritizes correctly identifying long-stay patients. This is particularly important in healthcare operations where **false positives for long stays can lead to unnecessary resource allocation**, whereas false negatives (missing short stays) are less critical in comparison. To achieve this, we fine-tuned our model using **Bayesian optimization**, focusing on precision as the key metric. The final model achieved a **precision of 0.89 for the long-stay class**, indicating that when the model predicts a long stay, it is correct 89% of the time. This high precision aligns well with our goal. While the model sacrifices recall—capturing fewer actual long-stay patients—this trade-off is acceptable within our context, as **minimizing false positives is more important than maximizing sensitivity**. Overall, the model is well-optimized for our use case and supports effective decision-making in predicting long hospital stays with high confidence.

## Future Scope

While the current model demonstrates strong precision for predicting long hospital stays, there are several promising directions to enhance and extend this work:

- 1. Improve Recall without Sacrificing Precision**  
Explore ensemble models or cost-sensitive learning techniques to improve recall for the long-stay class while maintaining high precision. This would help capture more true long-stay cases without increasing false positives.
- 2. Use Temporal Data from EHR**  
Integrate time-series data such as vital signs, lab results, and nursing notes over time using models like LSTMs or Transformers to better capture patient progression and risk of prolonged stay.
- 3. Leverage Multimodal Data**  
Combine structured data (e.g., lab values, vitals) with unstructured text (e.g., clinical notes) using multi-input models for a more holistic prediction pipeline.
- 4. Incorporate External Factors**  
Include social determinants of health, insurance type, and hospital resource availability, which often influence length of stay, to improve model generalizability.
- 5. Model Calibration & Explainability**  
Apply techniques like SHAP or LIME to better understand model predictions, improving clinician trust and adoption. Calibrated probability outputs can also aid in risk stratification.
- 6. Deploy in Real-time Settings**  
Integrate the model into clinical decision support systems for real-time prediction and alerts, enabling proactive discharge planning and resource allocation.
- 7. Cross-institutional Validation**  
Test the model on datasets from other hospitals to assess its robustness and adaptability to different patient populations and practices.