# Automating Data Cleaning for Healthcare Records using NLP

**Introduction:**

In modern healthcare systems, Electronic Health Records (EHRs) are used to store patient data digitally. However, these records often contain unstructured, inconsistent, or error-prone text—making analysis and modeling difficult. To improve the quality of healthcare data, this project applies Natural Language Processing (NLP) and Machine Learning (ML) techniques to automate the data cleaning process. The focus is on standardizing medical terms, correcting typographical errors, and ensuring consistency in patient records using AI-powered tools.

**Objectives:**

- To standardize medical terminology in EHRs.
- To detect and correct typographical errors in patient records.
- To ensure consistency across healthcare data using AI models.
- To build a pipeline using NLP tools that automates these tasks.

**Tools and packages required for project:**

**1. Jupyter Notebook:**

Jupyter Notebook is an open-source web application that allows users to create and share documents that contain live code, equations, visualizations, and narrative text. It is widely used in data science and machine learning for exploratory data analysis, modeling, and visualization.

Need of Jupyter notebook:

- To write and test Python code for NLP and ML tasks.
- To document the data cleaning steps interactively.
- To visualize outputs and track the performance of different models in one place.

**2. spaCy**

- spaCy is a fast and efficient Natural Language Processing (NLP) library in Python designed for real-world text processing. It supports tokenization, named entity recognition, and text normalization.
- Role in Project: It processes healthcare text, extracts important terms, and cleans the text using language rules, helping in standardizing medical records.

### 3. icd10

- The icd10 Python package provides access to the International Classification of Diseases (ICD-10) coding system used in healthcare. It allows lookup and validation of diagnosis codes.
- Role in Project: It maps raw diagnosis text to standardized ICD-10 codes, ensuring uniformity and aiding in healthcare data integration and analysis.

### 4. fuzzywuzzy

- fuzzywuzzy is a Python library that uses Levenshtein Distance to compare and match strings based on similarity. It's useful for handling small text variations.
- Role in Project: It identifies and corrects spelling mistakes and inconsistent terminology in medical records by matching similar terms.

### 5. dateutil

- dateutil is a powerful extension to Python's datetime module that helps parse, manipulate, and format dates from text.
- **Role in Project:** It cleans and standardizes dates in patient records, ensuring uniform formats for treatment history, admission dates, and more.

## Steps:

### 1. Data Collection

Collect raw Electronic Health Records (EHRs) in text format.

These may include patient notes, diagnosis summaries, and date entries.

### 2. Environment Setup

Install necessary tools and packages:

spaCy, icd10, fuzzywuzzy, dateutil, and Jupyter Notebook.

Load required NLP models (e.g., spaCy's English model).

### 3. Text Preprocessing

removing punctuation, special characters, and stopwords

Tokenizing and lemmatizing medical text using spaCy

### 4. Medical Term Standardization

Match terms to standard vocabulary:

Use fuzzywuzzy to correct typos and misspellings

Use icd10 to map diagnoses to standard ICD-10 codes

### 5. Date Normalization

Use dateutil to parse and standardize all date formats.

Ensure uniform date formats (e.g., DD-MM-YYYY) across records.

### 6. Result Visualization & Validation

Display the cleaned and standardized records.

Compare before-and-after data to check for improvements and correctness.

### 7. Documentation & Reporting

Document the process and outcomes in Jupyter Notebook.

Highlight examples, challenges faced, and final results.

Automating data cleaning for healthcare records using NLP improves the accuracy and efficiency of managing medical data. By extracting relevant information from unstructured sources and standardizing terms, it enhances data quality and ensures compliance with privacy regulations. This automation speeds up data processing, making it easier to analyze for better decision-making and patient outcomes. Ultimately, it reduces manual effort and supports scalable, secure healthcare operations.