

**K. J. Somaiya School of Engineering, Mumbai-77 TRUST**  
(Somaiya Vidyavihar University)  
**Department of Computer Engineering**

<b>Course Name:</b>	<b>Data Analysis Laboratory (216H03L501 )</b>	<b>Semester:</b>	<b>V</b>
<b>Date of Performance:</b>	<b>27 / 10 / 2025</b>	<b>DIV/ Batch No:</b>	<b>HDA2</b>
<b>Student Name:</b>	<b>Aaryan Sharma</b>	<b>Roll No:</b>	<b>16010123012</b>

**TITLE : NLP on clinical data (Finding data from Literature)**

**AIM:** Implement NLP on clinical data

**Expected Outcome of Experiment:**

**Books/ Journals/ Websites referred:**

**Sample case Study**

<https://towardsdatascience.com/clinical-named-entity-recognition-using-spacy-5ae9c002e86f>

**Theory:**

**Named entity recognition (NER)** is a natural language processing (NLP) method that extracts information from text. NER involves detecting and categorizing important information in text known as named entities. Named entities refer to the key subjects of a piece of text, such as names, locations, companies, events and products, as well as themes, topics, times, monetary values and percentages.

NER is also referred to as entity extraction, chunking and identification. It's used in many fields in artificial intelligence (AI), including machine learning (ML), deep learning and neural networks. NER is a key component of NLP systems, such as chatbots, sentiment analysis tools and search engines. It's used in healthcare, finance, human resources (HR), customer support, higher education and social media analysis.

**The purpose of NER**

NER identifies, categorizes and extracts the most important pieces of information from unstructured text without requiring time-consuming human analysis. It's particularly useful for quickly extracting key information from large amounts of data because it automates the extraction process.

As NER models improve their ability to correctly identify important information, they are helping improve AI systems in general. These systems are enhancing AI language comprehension capabilities in areas such as summarization and translation systems and the ability of AI systems to analyze text.

NER uses algorithms that function based on grammar, statistical NLP models and predictive models. These algorithms are trained on data sets that people label with predefined named entity categories, such as people, locations, organizations, expressions, percentages and monetary values. Categories are identified with

**K. J. Somaiya School of Engineering, Mumbai-77**  
(Somaiya Vidyavihar University)  
**Department of Computer Engineering**

abbreviations; for example, LOC is used for location, PER for persons and ORG for organizations.

**Dataset:**

<https://www.mtsamples.com>

<https://www.kaggle.com/datasets/tboyle10/medicaltranscriptions>

**Language: C/C++/Java/Python/other – Choice of student**

- Steps:**
1. Download dataset (Given).
  2. Pre-process the dataset if needed.
  3. Use spaCy library or any other similar one
  4. Visualize and Display result

**List out the library used with justification**

- spaCy: Used for performing Named Entity Recognition (NER) and text preprocessing. It provides pre-trained NLP pipelines and models that support biomedical and clinical domains through extensions like scispaCy.
- pandas: Used for reading, cleaning, and organizing text data
- matplotlib: Used for visualization of entity frequency or category distribution in the extracted data.

**Implementation details**

```
[22] ✓ 6s !pip install spacy

[23] ✓ 1m !pip install https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.5.0/en_core_web_sm-3.5.0.tar.gz
Collecting https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.5.0/en_core_web_sm-3.5.0.tar.gz
  Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.5.0/en_core_web_sm-3.5.0.tar.gz (12.8 MB)
    12.8/12.8 MB 39.7 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting spacy<3.6.0,>=3.5.0 (from en_core_web_sm==3.5.0)
  Using cached spacy-3.5.4.tar.gz (1.2 MB)

[29] ✓ 0s
# Sample healthcare-related text
healthcare_text = """
    The patient was diagnosed with type 2 cancer and malaria.
    They were prescribed benzodiazepines and ibuprofen. Further, the patient will undergo
    chemotherapy at Reliance Hospital next month."""

[30] ✓ 0s
# Process the text
doc = nlp(healthcare_text)
```

**K. J. Somaiya School of Engineering, Mumbai-77**  
(Somaiya Vidyavihar University)  
**Department of Computer Engineering**

```
[31]
✓ Os
# Print named entities
print("Entities in the text:")
for ent in doc.ents:
    print(f"{ent.text} - {ent.label_}")
```

Entities in the text:  
2 - CARDINAL  
Reliance Hospital - ORG  
next month - DATE

```
[32]
✓ Os
# Print explanations of entity labels
print("\nEntity Label Explanations:")
for ent in doc.ents:
    print(f"{ent.text}: {spacy.explain(ent.label_)}")
```

Entity Label Explanations:  
2: Numerals that do not fall under another type  
Reliance Hospital: Companies, agencies, institutions, etc.  
next month: Absolute or relative dates or periods

```
[38]
✓ Os
# Define custom entity labels for diseases and drugs
disease_labels = ["type 2 cancer", "malaria"] # Add more disease terms as needed
drug_labels = ["benzodiazepines", "ibuprofen"] # Add more drug names as needed
```

```
[39]
✓ Os
# Iterate through the text and identify entities based on the custom labels
entities = []
for label in disease_labels:
    if label in healthcare_text:
        entities.append((label, "DISEASE"))
for label in drug_labels:
    if label in healthcare_text:
        entities.append((label, "DRUG"))
```

```
[40]
✓ Os
# Print the identified entities
print("Entities in the text:")
for entity, label in entities:
    print(f"{entity} - {label}")
```

Entities in the text:  
type 2 cancer - DISEASE  
malaria - DISEASE  
benzodiazepines - DRUG  
ibuprofen - DRUG

**K. J. Somaiya School of Engineering, Mumbai-77** TRUST  
(Somaiya Vidyavihar University)  
**Department of Computer Engineering**

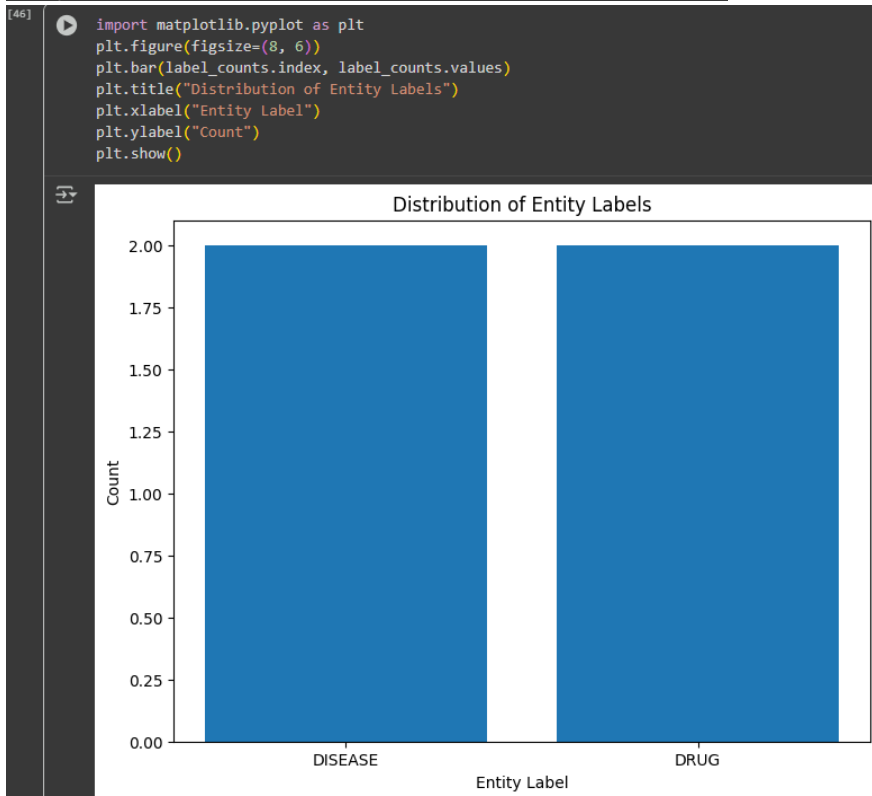
```
[44] import pandas as pd
entity_df = pd.DataFrame(entities, columns=['Entity', 'Label'])
display(entity_df)
```

	Entity	Label
0	type 2 cancer	DISEASE
1	malaria	DISEASE
2	benzodiazepines	DRUG
3	ibuprofen	DRUG

```
[45] label_counts = entity_df['Label'].value_counts()
display(label_counts)
```

	count
DISEASE	2
DRUG	2

dtype: int64



**K. J. Somaiya School of Engineering, Mumbai-77** TRUST  
(Somaiya Vidyavihar University)  
**Department of Computer Engineering**

**Conclusion (Interpretation of result):**

The experiment successfully demonstrates the use of Named Entity Recognition (NER) for extracting clinical entities from unstructured medical text. Using spaCy, we efficiently identified medical terms such as diseases, drugs, and treatments from clinical notes. This automation reduces manual data extraction effort and provides structured insights from large volumes of clinical documents. Such NLP-based techniques are highly valuable in clinical data analysis, medical research, and healthcare informatics.