

Data Science Project

Dr. Mohammad Al-Hammouri



Predicting ICD-9 Diagnosis Codes from Discharge Summaries using NLP

Moath Yousef AL-Tahrawi	2141099
Ahmad Yadak	2135721
Malik Alarja	2141518

ABSTRACT

For phase one, we were provided with four (.csv) data sets consisting of: ADMISSIONS, D_ICD_DIAGNOSSES, DIAGNOSSES_ICD_SORTED, PATIENTS. and asked to preprocess them, by plotting the data columns, dropping the unwanted columns, binarizing the needed data, handling the empty slots, etc.

Beginning with merging the files together, we combined all four data sets depending on their foreign key and primary key connections, resulting in the final file named MERGED_PROCESSED_DATA_FINAL, then we preprocessed each column depending on what we think is needed or not, with taking into consideration the plotting of the related columns.

In phase two, the project focuses on predicting ICD-9 diagnosis codes from hospital discharge summaries using Natural Language Processing (NLP). By leveraging clinical text data from the MIMIC dataset, we aim to develop models that can automatically assign one or more ICD-9 codes to each discharge note. The workflow involves preprocessing the text, aligning each note with its corresponding diagnosis codes, and training classification models. We compare the performance of transformer-based models like ClinicalBERT with deep learning models such as CNNs and LSTMs. Evaluation is based on standard multi-label classification metrics, including precision, recall, F1-score, Hamming loss, and top K accuracy. This work supports the automation of clinical coding, which is vital for hospital billing, analytics, and improving patient care efficiency.

PHASE 1: Data Cleaning and Preprocessing

As mentioned before, in this part of the project we worked on cleaning the four data sets:

- 1) ADMISSIONS
- 2) D_ICD_DIAGNOSSES
- 3) DIAGNOSSES_ICD_SORTED
- 4) PATIENTS

Part A: Merging the data sets

- 1) Load each data file (e.g., CSV format) into a pandas DataFrame.

```
import pandas as pd

# Load the CSV files
admissions_df = pd.read_csv("/content/ADMISSIONS.csv")
icd_diagnoses_df = pd.read_csv("/content/D_ICD_DIAGNOSSES.csv")
diagnoses_icd_df = pd.read_csv("/content/DIAGNOSSES_ICD_sorted.csv")
patients_df = pd.read_csv("/content/PATIENTS.csv")
```

- 2) Merge diagnoses with ICD descriptions, then merge with admissions and patient data.

```
# Step 1: Merge diagnoses with ICD descriptions
merged_df = diagnoses_icd_df.merge(icd_diagnoses_df, on="ICD9_CODE", how="left")

# Step 2: Merge with admissions data
merged_df = merged_df.merge(admissions_df, on=["SUBJECT_ID", "HADM_ID"], how="left")

# Step 3: Drop duplicate or irrelevant ROW_ID columns before merging
if 'ROW_ID' in patients_df.columns:
    patients_df = patients_df.drop(columns=['ROW_ID'])
if 'ROW_ID' in merged_df.columns:
    merged_df = merged_df.drop(columns=['ROW_ID'])

# Now merge with patient demographics
merged_df = merged_df.merge(patients_df, on="SUBJECT_ID", how="left")
```

Note: Before merging PATIENTS, duplicates are removed to **avoid duplicate/conflicting column names** in the final merged dataset.

Result of part A:

```
Dropped columns: ['SUBJECT_ID', 'ROW_ID_x']
  HADM_ID  SEQ_NUM  ICD9_CODE  ROW_ID_y  SHORT_TITLE \
0  163353      1.0      V3001  13695.0  Single lb in-hosp w cs
1  163353      2.0      V053  12202.0  Need prphyl vc vrl hepat
2  163353      3.0      V290  13688.0  NB obsrv suspct infect
3  145834      9.0      2639   1642.0  Protein-cal malnutr NOS
4  145834      7.0      6826   7283.0  Cellulitis of leg

  LONG_TITLE  ADMITTIME \
0  Single liveborn, born in hospital, delivered b... 2138-07-17 19:04:00
1  Need for prophylactic vaccination and inoculat... 2138-07-17 19:04:00
2  Observation for suspected infectious condition 2138-07-17 19:04:00
3  Unspecified protein-calorie malnutrition 2101-10-20 19:08:00
4  Cellulitis and abscess of leg, except foot 2101-10-20 19:08:00

  DISCHTIME  DEATHTIME  ADMISSION_TYPE  ...  DIAGNOSIS \
0  2138-07-21 15:48:00      NaN      NEWBORN  ...  NEWBORN
1  2138-07-21 15:48:00      NaN      NEWBORN  ...  NEWBORN
2  2138-07-21 15:48:00      NaN      NEWBORN  ...  NEWBORN
3  2101-10-31 13:58:00      NaN      EMERGENCY  ...  HYPOTENSION
4  2101-10-31 13:58:00      NaN      EMERGENCY  ...  HYPOTENSION

  HOSPITAL_EXPIRE_FLAG  HAS_CHARTEVENTS_DATA  GENDER  DOB  DOD \
0      0      1  M  2138-07-17      NaN
1      0      1  M  2138-07-17      NaN
2      0      1  M  2138-07-17      NaN
3      0      1  M  2025-04-11  2102-06-14
4      0      1  M  2025-04-11  2102-06-14

  DOD_HOSP  DOD_SSN  EXPIRE_FLAG  AGE
0      NaN      NaN      0      0
1      NaN      NaN      0      0
2      NaN      NaN      0      0
3      NaN  2102-06-14      1      76
4      NaN  2102-06-14      1      76
[5 rows x 29 columns]
```

PART B: Data Preprocessing

- 1) Convert relevant date columns to datetime objects.

```
# Convert date columns to datetime
date_columns = [
    "ADMITTIME", "DISCHTIME", "DEATHTIME",
    "DOB", "DOD", "DOD_HOSP", "DOD_SSN",
    "EDREGTIME", "EDOUTTIME"
]
for col in date_columns:
    if col in merged_df.columns:
        merged_df[col] = pd.to_datetime(merged_df[col], errors='coerce')
```

- 2) Fill missing categorical values with "Unknown".

```
# Fill missing categorical data with 'Unknown'
categorical_cols = [
    "ADMISSION_TYPE", "ADMISSION_LOCATION", "DISCHARGE_LOCATION",
    "INSURANCE", "LANGUAGE", "RELIGION", "MARITAL_STATUS",
    "ETHNICITY", "SHORT_TITLE", "LONG_TITLE", "DIAGNOSIS"
]
for col in categorical_cols:
    if col in merged_df.columns:
        merged_df[col] = merged_df[col].fillna("Unknown")
```

- 3) Fill missing numeric values (e.g., sequence numbers) appropriately and Reset DataFrame index and save intermediate merged data.

```
# Fill numeric NaNs like SEQ_NUM
if 'SEQ_NUM' in merged_df.columns:
    merged_df['SEQ_NUM'] = merged_df['SEQ_NUM'].fillna(0)

merged_df.reset_index(drop=True, inplace=True)

merged_df.to_csv("/content/merged_preprocessed_data.csv", index=False)
```

Result of part B:

```
Added LENGTH_OF_STAY feature:
  ADMITTIME DISCHTIME LENGTH_OF_STAY
0 2138-07-17 19:04:00 2138-07-21 15:48:00      3.0
1 2138-07-17 19:04:00 2138-07-21 15:48:00      3.0
2 2138-07-17 19:04:00 2138-07-21 15:48:00      3.0
3 2138-07-17 19:04:00 2138-07-21 15:48:00      3.0
4 2138-07-17 19:04:00 2138-07-21 15:48:00      3.0

Scaled Numerical Features:
  AGE LENGTH_OF_STAY
0 -2.460285 -0.665227
1 -2.460285 -0.665227
2 -2.460285 -0.665227
3 -2.460285 -0.665227
4 -2.460285 -0.665227

DataFrame after One-Hot Encoding:
  HADM_ID  SEQ_NUM  ICD9_CODE  ROM_ID_y  SHORT_TITLE \
0  163353      1.0    V3001  13695.0  Single lb in-hosp w cs
1  163353      1.0    V3001  13695.0  Single lb in-hosp w cs
2  163353      1.0    V3001  13695.0  Single lb in-hosp w cs
3  163353      2.0    V053   12202.0  Need prphyl vc vrl hepat
4  163353      2.0    V053   12202.0  Need prphyl vc vrl hepat

  LONG_TITLE  ADMITTIME \
0 Single liveborn, born in hospital, delivered b... 2138-07-17 19:04:00
1 Single liveborn, born in hospital, delivered b... 2138-07-17 19:04:00
2 Single liveborn, born in hospital, delivered b... 2138-07-17 19:04:00
3 Need for prophylactic vaccination and inoculat... 2138-07-17 19:04:00
4 Need for prophylactic vaccination and inoculat... 2138-07-17 19:04:00

  DISCHTIME DEATHTIME  ADMISSION_LOCATION  ...  AGE \
0 2138-07-21 15:48:00  NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
1 2138-07-21 15:48:00  NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
2 2138-07-21 15:48:00  NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
3 2138-07-21 15:48:00  NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
4 2138-07-21 15:48:00  NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285

  LENGTH_OF_STAY  GENDER_M  ADMISSION_TYPE_EMERGENCY  ADMISSION_TYPE_NEWBORN \
0 -0.665227      True      False      True
1 -0.665227      True      False      True
2 -0.665227      True      False      True
3 -0.665227      True      False      True
4 -0.665227      True      False      True

  ADMISSION_TYPE_URGENT  INSURANCE_Medicaid  INSURANCE_Medicare \
0      False      False      False
1      False      False      False
2      False      False      False
3      False      False      False
4      False      False      False

  INSURANCE_Private  INSURANCE_Self Pay
0      True      False
1      True      False
2      True      False
3      True      False
4      True      False
```

PART C: Calculate Patient Age and Plot Distributions

```
# =====  
# PART 3: Calculate Age and Filter  
# =====  
  
import matplotlib.pyplot as plt  
import seaborn as sns  
  
# Reload if needed  
merged_df = pd.read_csv("/content/merged_preprocessed_data.csv")  
  
merged_df["ADMITTIME"] = pd.to_datetime(merged_df["ADMITTIME"], errors="coerce")  
merged_df["DOB"] = pd.to_datetime(merged_df["DOB"], errors="coerce")  
  
merged_df = merged_df.dropna(subset=["DOB", "ADMITTIME"])  
merged_df = merged_df[merged_df["DOB"] < merged_df["ADMITTIME"]]  
  
year_diff = merged_df["ADMITTIME"].dt.year - merged_df["DOB"].dt.year  
age = year_diff - ((merged_df["ADMITTIME"].dt.month < merged_df["DOB"].dt.month) |  
                  ((merged_df["ADMITTIME"].dt.month == merged_df["DOB"].dt.month) &  
                   (merged_df["ADMITTIME"].dt.day < merged_df["DOB"].dt.day)))  
merged_df["AGE"] = age  
  
merged_df = merged_df[(merged_df["AGE"] >= 0) & (merged_df["AGE"] <= 120)]
```

PART D: Plot Age and Categorical Distributions

```
# Plot Age histogram  
plt.figure(figsize=(10, 5))  
sns.histplot(merged_df["AGE"], bins=30, kde=True, color='skyblue')  
plt.title("Distribution of Patient Age at Admission")  
plt.xlabel("Age")  
plt.ylabel("Number of Patients")  
plt.grid(True)  
plt.show()  
  
# Plot Age violin plot by gender  
plt.figure(figsize=(8, 5))  
sns.violinplot(x="GENDER", y="AGE", data=merged_df, palette="pastel", inner="box")  
plt.title("Violin Plot of Age Distribution by Gender")  
plt.ylabel("Age")  
plt.xlabel("Gender")  
plt.grid(True)  
plt.show()
```

PART E: Drop Unwanted Columns

```
columns_to_drop = ['SUBJECT_ID', 'ROW_ID_X']
existing_to_drop = [col for col in columns_to_drop if col in merged_df.columns]

if existing_to_drop:
    merged_df = merged_df.drop(columns=existing_to_drop, errors='ignore')
    print(f"Dropped columns: {existing_to_drop}")
else:
    print("No specified columns to drop found.")

print(merged_df.head())
```

PART F: Feature Engineering, Scaling, Encoding

1) Calculate Length of Stay

```
# Calculate Length of Stay
if 'ADMITTIME' in merged_df.columns and 'DISCHTIME' in merged_df.columns:
    valid_los_df = merged_df[merged_df['DISCHTIME'] >= merged_df['ADMITTIME']].copy()
    valid_los_df['DISCHTIME'] = pd.to_datetime(valid_los_df['DISCHTIME'], errors='coerce')
    valid_los_df['ADMITTIME'] = pd.to_datetime(valid_los_df['ADMITTIME'], errors='coerce')
    valid_los_df['LENGTH_OF_STAY'] = (valid_los_df['DISCHTIME'] - valid_los_df['ADMITTIME']).dt.days
    merged_df = merged_df.merge(valid_los_df[['HADM_ID', 'LENGTH_OF_STAY']], on='HADM_ID', how='left')
    merged_df['LENGTH_OF_STAY'] = merged_df['LENGTH_OF_STAY'].fillna(0)
    print("\nAdded LENGTH_OF_STAY feature:")
    print(merged_df[['ADMITTIME', 'DISCHTIME', 'LENGTH_OF_STAY']].head())
```

2) Scale numerical features

```
from sklearn.preprocessing import StandardScaler

num_features = ['AGE', 'LENGTH_OF_STAY']
num_features = [f for f in num_features if f in merged_df.columns]

if num_features:
    scaler = StandardScaler()
    merged_df[num_features] = scaler.fit_transform(merged_df[num_features])
    print("\nScaled Numerical Features:")
    print(merged_df[num_features].head())
else:
    print("\nNumerical features not found for scaling.")
```

3) One-hot encode categorical variables

```
merged_df = pd.get_dummies(merged_df, columns=['GENDER', 'ADMISSION_TYPE', 'INSURANCE'], drop_first=True)
print("\nDataFrame after One-Hot Encoding:")
print(merged_df.head())
print(merged_df.columns)

merged_df.reset_index(drop=True, inplace=True)
merged_df.to_csv("/content/merged_preprocessed_data_final.csv", index=False)
```


PART G: Convert Boolean Columns to Integers and Save

```
merged_df = merged_df.replace({True: 1, False: 0})

print("\nDataFrame after converting True/False to 1/0:")
print(merged_df.head())

merged_df.to_csv("/content/merged_preprocessed_data_final.csv", index=False)
```

Result of part G :

```
DataFrame after converting True/False to 1/0:
  HADM_ID  SEQ_NUM  ICD9_CODE  ROW_ID_y  SHORT_TITLE \
0  163353    1.0    V3001    13695.0    Single lb in-hosp w cs
1  163353    1.0    V3001    13695.0    Single lb in-hosp w cs
2  163353    1.0    V3001    13695.0    Single lb in-hosp w cs
3  163353    2.0    V053    12202.0    Need prphyl vc vr1 hepat
4  163353    2.0    V053    12202.0    Need prphyl vc vr1 hepat

  LONG_TITLE  ADMITTIME \
0  Single liveborn, born in hospital, delivered b...  2138-07-17 19:04:00
1  Single liveborn, born in hospital, delivered b...  2138-07-17 19:04:00
2  Single liveborn, born in hospital, delivered b...  2138-07-17 19:04:00
3  Need for prophylactic vaccination and inoculat...  2138-07-17 19:04:00
4  Need for prophylactic vaccination and inoculat...  2138-07-17 19:04:00

  DISCHTIME  DEATHTIME  ADMISSION_LOCATION  ...  AGE \
0  2138-07-21 15:48:00    NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
1  2138-07-21 15:48:00    NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
2  2138-07-21 15:48:00    NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
3  2138-07-21 15:48:00    NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285
4  2138-07-21 15:48:00    NaN  PHYS REFERRAL/NORMAL DELI  ... -2.460285

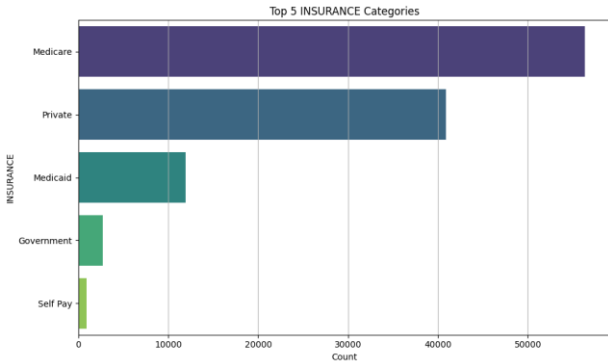
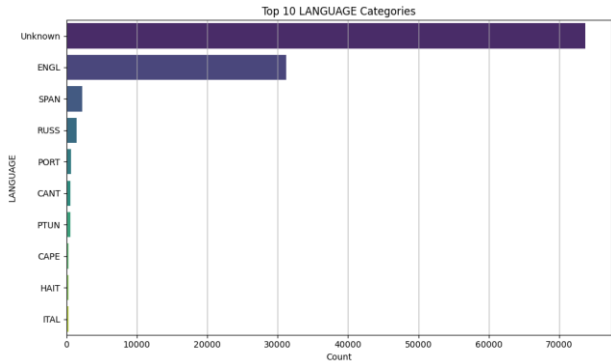
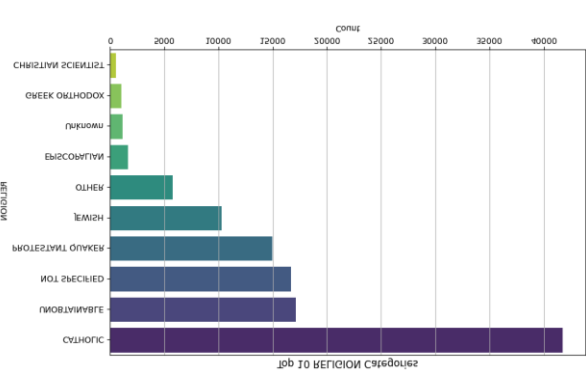
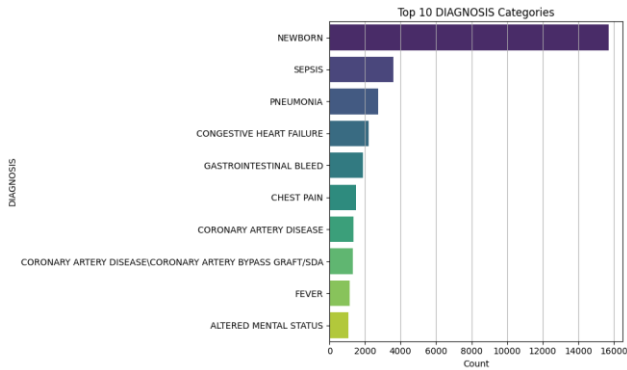
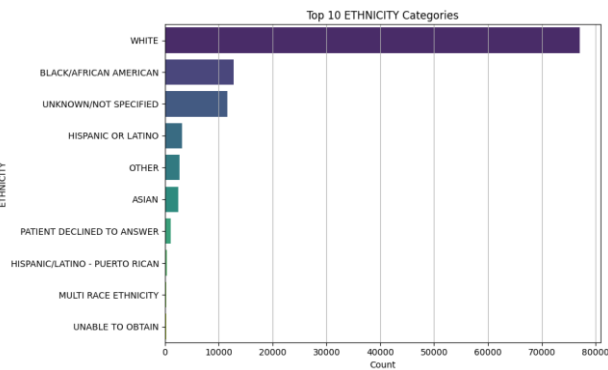
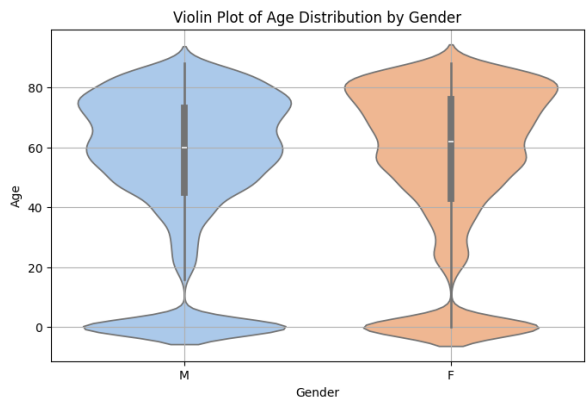
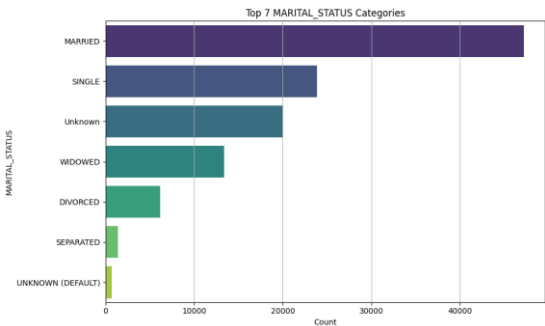
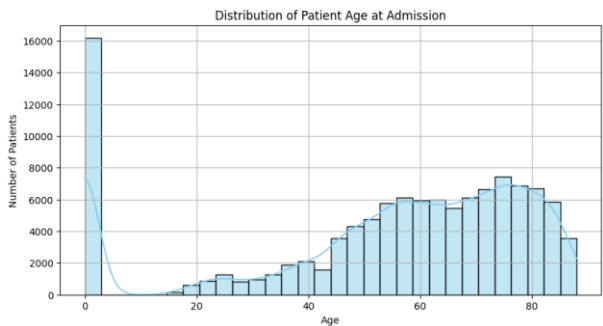
  LENGTH_OF_STAY  GENDER_M  ADMISSION_TYPE_EMERGENCY  ADMISSION_TYPE_NEWBORN \
0      -0.665227         1              0              1
1      -0.665227         1              0              1
2      -0.665227         1              0              1
3      -0.665227         1              0              1
4      -0.665227         1              0              1

  ADMISSION_TYPE_URGENT  INSURANCE_Medicaid  INSURANCE_Medicare \
0              0              0              0
1              0              0              0
2              0              0              0
3              0              0              0
4              0              0              0

  INSURANCE_Private  INSURANCE_Self Pay
0              1              0
1              1              0
2              1              0
3              1              0
4              1              0
```


PHASE ONE RESULTS:

PLOTS GENERATED:



PHASE 2: Predicting ICD-9 Diagnosis Codes from Discharge Summaries using NLP

The task is to train two different models on a certain data set

MODELS :

- 1) CLINICAL-BERT : a domain-specific variant of **BERT (Bidirectional Encoder Representations from Transformers)** that has been further **pretrained on clinical text**, specifically from the **MIMIC-III** dataset.
- 2) LSTM (LONG SHORT-TERM MEMORY) : is a type of **Recurrent Neural Network (RNN)** that is especially good at learning from **sequential data**, such as text. When working with a **clinical** dataset, LSTMs can help model the **sequence of words in clinical notes** (like discharge summaries) to predict **ICD-9 codes**.

DATA SET :

The dataset “rntc/mimic-icd-visit” comes from the **MIMIC-III** clinical database and is provided through the Hugging Face Datasets library by the **Research on NLP for Text Classification (RNTC)** project.

- 1) It contains data from over **40,000 ICU patients** from the **Beth Israel Deaconess Medical Center** between 2001–2012.
- 2) The dataset includes demographics, vital signs, medications, lab measurements, and **clinical notes** (such as discharge summaries).

Each entry in the dataset represents a **hospital visit** and includes:

subject_id: Unique ID for each patient

hadm_id: Hospital admission ID for the specific visit

cleaned_text: A cleaned version of the **discharge summary** note

icd_code: A list of **ICD-9 diagnosis codes** assigned for that visit

NOTE: WE ARE GOING TO USE THE ICD_CODE AND CLEANED_TEXT ATTRIBUTES IN THIS TASK TO PREDICT THE ICD_CODE FROM THE CLEANED TEXT(DISCHARGE SUMMARY IN OTHER WORDS).

FIRST MODEL: Clinical Bert

ClinicalBERT is a domain-specific version of the BERT language model that is fine-tuned to understand clinical and medical text. Standard BERT is pretrained on general-domain corpora like Wikipedia and books. While powerful, it **lacks understanding of medical jargon, abbreviations, and syntax** found in clinical notes.

ClinicalBERT bridges this gap by **continuing BERT's pretraining on clinical notes** from the **MIMIC-III** database.

This makes ClinicalBERT highly specialized in:

- Medical terminology
- Clinical grammar and abbreviations
- Context in healthcare narratives

You can load ClinicalBERT from Hugging Face:

```
from transformers import AutoTokenizer, AutoModel

tokenizer = AutoTokenizer.from_pretrained("emilyalsentzer/Bio_ClinicalBERT")
model = AutoModel.from_pretrained("emilyalsentzer/Bio_ClinicalBERT")
```

Results:

```
Epoch 15/15
Training: 100%|██████████| 2813/2813 [06:49<00:00, 6.86it/s]
Validating: 100%|██████████| 313/313 [00:40<00:00, 7.65it/s]
Train - loss: 0.0044, acc: 0.9988, f1: 0.2981, ham: 0.0000
Val   - loss: 0.0050, acc: 0.9987, f1: 0.2824, ham: 0.0000
```

Conclusion :

ClinicalBERT is a transformer-based language model tailored for healthcare, trained on real ICU patient notes from MIMIC-III. By capturing the structure and terminology of clinical text, it outperforms general BERT models in medical NLP tasks. It is widely used for automating clinical coding, understanding discharge summaries, and improving decision support in healthcare systems.

SECOND MODEL: LONG SHORT-TERM MEMORY

LSTM (Long Short-Term Memory) is a type of Recurrent Neural Network (RNN) that is designed to learn from sequential data — such as text — by remembering long-term dependencies. When applied to clinical data like discharge summaries from the mimic-icd-visit dataset, LSTM can learn how patterns in medical language relate to ICD-9 diagnosis codes.

- 1) Clinical notes are **long, unstructured, and sequential** in nature.
- 2) LSTM can process text **word by word**, preserving **context over long distances**.
- 3) It avoids issues like the **vanishing gradient problem** seen in standard RNN

Key Points

- LSTMs are ideal for learning long-term patterns in clinical text.
- Unlike BERT, LSTM **does not require transformer-level compute resources**.
- Can be enhanced with **BiLSTM (Bidirectional LSTM)** to capture both past and future context.
- Can be stacked or combined with **CNN layers** or **attention mechanisms** for better performance.

Results:

```
Epoch 9 Training: 100%|██████████| 424/424 [05:47<00:00, 1.22it/s]
Epoch 9 - Loss: 0.0088 - F1 Micro: 0.2156, Macro: 0.0056, average: 0.1095 -
Recall Micro: 0.2779, Macro: 0.0146, average: 0.2779 -
Precision Micro: 0.1761, Macro: 0.0038, average: 0.0721 - Top-10 Accuracy: 0.8627
```

Conclusion

LSTM is an effective and efficient deep learning architecture for modeling clinical text in the mimic-icd-visit dataset. Its ability to retain contextual information across long discharge summaries makes it well-suited for predicting multiple ICD-9 codes. While newer models like BERT outperform in many cases, LSTM remains a strong and interpretable baseline for sequence-based medical NLP tasks.

Theoretical comparison:

Feature	ClinicalBERT	LSTM Model
Architecture	Transformer-based (BERT pretrained on clinical notes)	Recurrent Neural Network (LSTM)
Pretraining	Yes – trained on biomedical/caregiver notes	No – typically trained from scratch
Multi-label Capability	Strong – captures dependencies between labels well	Moderate – less effective for complex or rare code combinations
Performance (F1 Score)	Higher – better contextual understanding and accuracy	Lower – especially struggles with long or complex sequences
Training Time	Longer – due to model size and complexity	Shorter – fewer parameters
Hardware Requirement	High – requires GPU for efficient training and inference	Moderate – can train on CPU for small-scale datasets
Fine-tuning	Yes – pretrained model needs fine-tuning on ICD-9 classification task	Yes – trained from scratch or fine-tuned if pretrained embeddings used
Best For	High-accuracy ICD-9 code prediction, complex clinical language	Lightweight setups, when resources are limited

Implementation Comparison:

- 1) In our models, they both got near f1 scores
- 2) The clinical bert model trained more dataset samples = 50000, while the LSTM was trained on 17000, because of hardware performance limitations we had in implementation, if we exceeded this number of samples the computer crashes, or takes many more time than needed to train.
- 3) The LSTM took less time in training (almost 6 minutes per epoch) while the clinical bert took more time (almost 16 minutes per epoch)
- 4) The input/output testing gave great results in both models
- 5) Both models were easy to edit and fine tuning feasible

ADD ONS:

We wrote a code (icd_checker_testing) that searches in the data set for a certain icd code, and returns the corresponding “Cleaned_text” (discharge submission) to implement multiple input/output testing scenarios.

```
# Function to retrieve cleaned_text for a given icd_code
def get_texts_for_icd_code(target_icd):
    results = []
    for example in ds:
        if target_icd in example["icd_code"]: # icd_code is a list
            results.append(example["cleaned_text"])
    return results

# Example usage
icd_input = "E785"
matched_texts = get_texts_for_icd_code(icd_input)

# Print first few matches
for i, text in enumerate(matched_texts[:5]):
    print(f"Example {i+1}: \n{text} \n{'-'*40}")
```

This part searches for the icd code (icd input variable)

```
for i in range(10):
    print(f"ICD Codes {i+1}: {ds[i]['icd_code']}")
```

Here you generate multiple icd codes to search for in the data set

The end