

# MIMIC NLP

## Tools and Dataset Used:

- Filename: [https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/mimic\\_nlp.ipynb](https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/mimic_nlp.ipynb)
- Dataset Used: [https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/nlp\\_med\\_notes.csv](https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/nlp_med_notes.csv)
- Package Used: Spacy, Scispacy, BioClinicalBERT, Word2Vec, TSNE
- Application Used: Dbeaver, VSCode

# Dataset Preparation:

- Write a sql query in Dbeaver
- Use the 'Export data' option in Dbeaver to extract the dataset result as CSV file into NLP workspace

```

89 with cte as (
90   select * from d_icd_diagnoses did where did.icd9_code between '25000' and '25003'),
91   cte1 as (select di.* from diagnoses_icd di inner join cte on cte.icd9_code = di.icd9_code ),
92   cte2 as (select cte1.icd9_code,n.* from noteevents n inner join cte1 on n.hadm_id = cte1.hadm_id ),
93   cte3 as (
94     select icd9_code,subject_id, hadm_id,category,text from cte2 where category Like'Discharge%' limit 40)
95   select row_number( ) over (order by hadm_id ) as rn, cte3.* from cte3;
96

```

diagnoses\_icd(+) 1 ×

with cte as (select \* from d\_icd\_diagnoses did where did.icd9\_code between '25000' and '25003')

	rn	icd9_code	subject_id	hadm_id	category	text
1	1	25000	157	107,880	Discharge summary	Admission Date: [**2106-6-17**] Discharge Date: [**2106-6-24**]¶¶Date of
2	2	25000	305	108,015	Discharge summary	Admission Date: [**2125-12-31**] Discharge Date: [**2126-1-10**]¶¶Date of
3	3	25000	21	109,451	Discharge summary	Admission Date: [**2134-9-11**] Discharge Date: [**2134-9-24**]¶¶Service:
4	4	25000	21	111,970	Discharge summary	Admission Date: [**2135-1-30**] Discharge Date: [**2135-2-8**]¶¶Service: h
5	5	25000	75	112,086	Discharge summary	Admission Date: [**2147-4-5**] Discharge Date: [**2147-4-11**]¶¶Date of Bi
6	6	25000	130	113,323	Discharge summary	Unit No: [**Numeric Identifier 56787**]¶¶Admission Date: [**2119-11-14**]¶¶Discharge I
7	7	25000	249	116,935	Discharge summary	Admission Date: [**2149-12-17**] Discharge Date: [**2149-12-31**]¶¶Date c
8	8	25000	188	132,401	Discharge summary	Admission Date: [**2161-11-1**] Discharge Date: [**2162-1-17**]¶¶Date of f
9	9	25000	305	133,059	Discharge summary	Admission Date: [**2125-4-26**] Discharge Date: [**2125-5-3**]¶¶Date of Bi
10	10	25000	205	135,671	Discharge summary	Admission Date: [**2191-11-6**] Discharge Date: [**2191-11-14**]¶¶Date of Birth
11	11	25000	205	135,671	Discharge summary	Admission Date: [**2191-11-6**] Discharge Date: [**2191-11-16**]¶¶Date of Birth
12	12	25000	191	136,614	Discharge summary	Admission Date: [**2196-4-9**] Discharge Date: [**2196-4-21**]¶¶Date of Bi
13	13	25000	184	137,477	Discharge summary	Admission Date: [**2168-3-13**] Discharge Date: [**2168-3-16**]¶¶Date of Birth:
14	14	25000	117	140,784	Discharge summary	Admission Date: [**2133-4-7**] Discharge Date: [**2133-4-12**]¶¶Date of Birth: ['
15	15	25000	13	143,045	Discharge summary	Name: [**Known lastname 9900**], [**Known firstname **] C Unit No: [**
16	16	25000	13	143,045	Discharge summary	Admission Date: [**2167-1-8**] Discharge Date: [**2167-1-15**]¶¶Date of Birth:
17	17	25000	249	149,546	Discharge summary	Admission Date: [**2155-2-3**] Discharge Date: [**2155-2-14**]¶¶Date of Bi
18	18	25000	294	152,578	Discharge summary	Admission Date: [**2118-1-17**] Discharge Date: [**2118-2-2**]¶¶Date of Bi

- The disease of interest here is diabetes
- The icd9\_code for diabetes is between '25000' and '25003'
- First we need to join d\_icd\_diagnoses and diagnoses\_icd
- The result of previous step subquery is joined with noteevents on hospital admission id.
- The notes of 'Discharge Summary' needs to be filtered out limited to 40 records in order to process effectively.

The screenshot shows the DBeaver 24.1.5 interface. The main window displays a SQL script in the 'Script Editor' tab. The script is a complex query involving multiple tables and joins, including 'Readmissions', 'caregivers', 'procedureevents\_mv', 'diagnoses\_icd', 'd\_icd\_diagnoses', and 'noteevents'. The script includes comments and SQL statements for selecting data, altering tables, and adding constraints.

A 'Data Transfer' dialog box is open in the foreground, showing the 'Export target' configuration. The 'Export target' is set to 'Database table(s)'. The 'Export settings' are configured for 'CSV' format. The 'Output' is set to 'Export to CSV file(s)'. The 'Confirm' button is highlighted.

The background window shows a table of data with columns for 'id', 'icd9\_code', 'admission\_date', 'discharge\_date', and 'note'. The table contains 40 rows of data, with the first row showing '11', '25000', '203', and '136,614'. The table is filtered to show 40 rows, with the first row highlighted.

# Named Entity Recognition:

- Import the Spacy, Scispacy, Word2Vec, TSNE libraries

```
import pandas as pd
pd.options.mode.chained_assignment = None
import numpy as np
import re
from gensim.models import Word2Vec
import gensim.downloader as api
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
from spacy import displacy
import spacy
spacy.require_gpu()
```

Python

- Load the spacy model

```
nlp = spacy.load('en_core_web_sm')
```

- Prepare utility functions like,
  - clean\_and\_split\_paragraph()
  - extract\_entities()
  - fetch\_entities()
  - visualize\_entities()
  - extract\_corpus()
  - fetch\_corpus()
  - extract\_passage\_by\_label()
  - do\_data\_process()

## Utility Helper Functions:

- **clean\_and\_split\_paragraph()** is used to perform data cleaning by removing the extra spaces and redundant lines.
- **fetch\_entities()** and **extract\_entities()** are used to print entities fetched using NER models.
- **fetch\_corpus()** and **extract\_corpus()** are used to print corpus recognized from the dataset.
- **do\_data\_process()** and **extract\_passage\_by\_label()** are used to extract subparagraph from medical notes passage.
- **visualize\_entities()** is used to identify the different types of entities provided by different **Spacy** and **SciSpacy** models
- For instance , **History of Present Illness, Past Medical History, Brief Hospital Course, REVIEW OF SYSTEMS** are the subparagraph labels

- Load and clean the 'nlp\_med\_notes.csv' dataset

```
notes_df = pd.read_csv("nlp_med_notes.csv")['text']
notes=[]
for notes_data in notes_df:
    notes.append(clean_and_split_paragraph(notes_data))
```

- For each record extract the 'PAST MEDICAL HISTORY' subparagraph from text (i.e medical notes) field.
- Visualize the extracted label\_df to view the entities of spacy model.

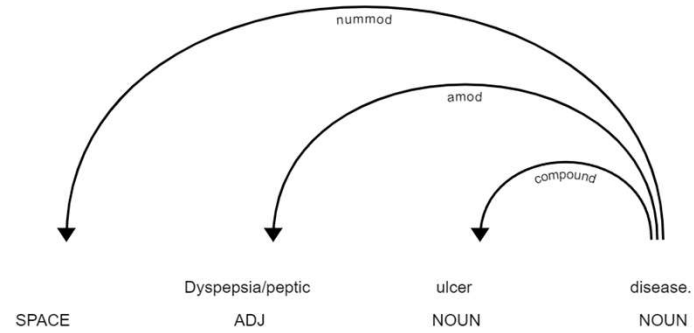
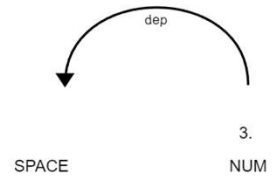
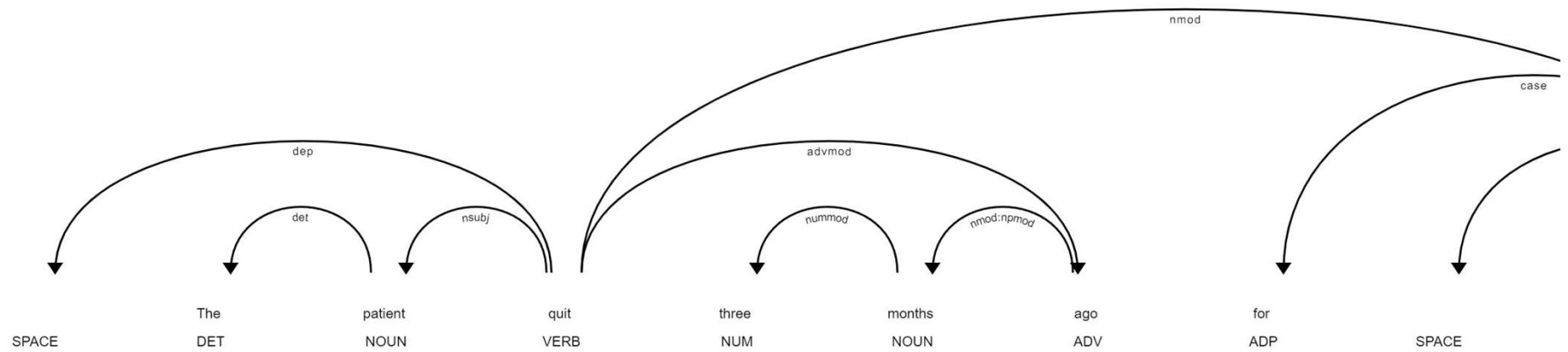
Hepatitis C genotype 1A ENTITY complicated ENTITY by  
 cirrhosis, esophageal varices ENTITY, encephalopathy ENTITY and  
 presumptive SBP ENTITY. Type 2 diabetes. Obesity ENTITY.  
 Hypertension ENTITY. Asthma.  
 Esophageal candidiasis ENTITY. Gastroparesis. Depression. Status  
 post cholecystectomy ENTITY. Status post ENTITY seven spur  
 surgeries ENTITY.  
 Hypothyroidism ENTITY. Amenorrhea. Migraines ENTITY.

\*\*\*\*\*

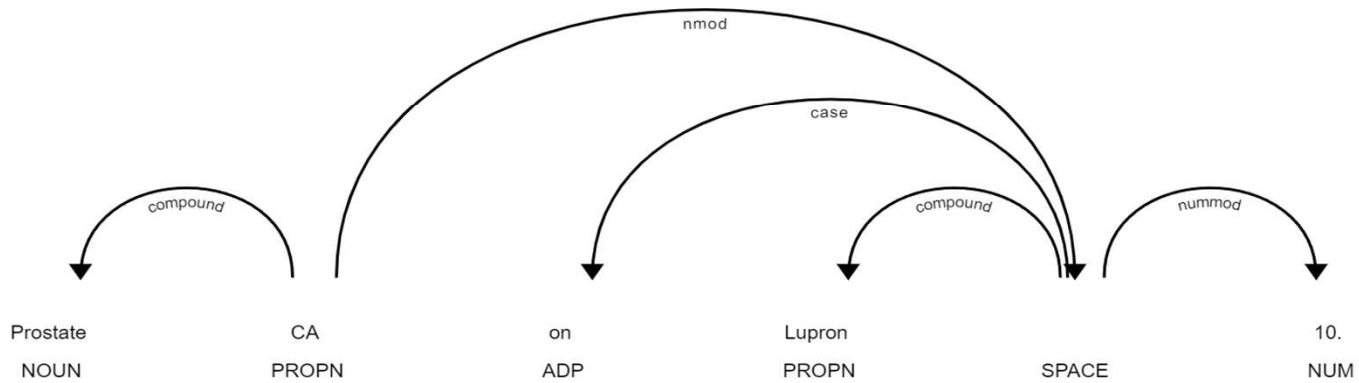
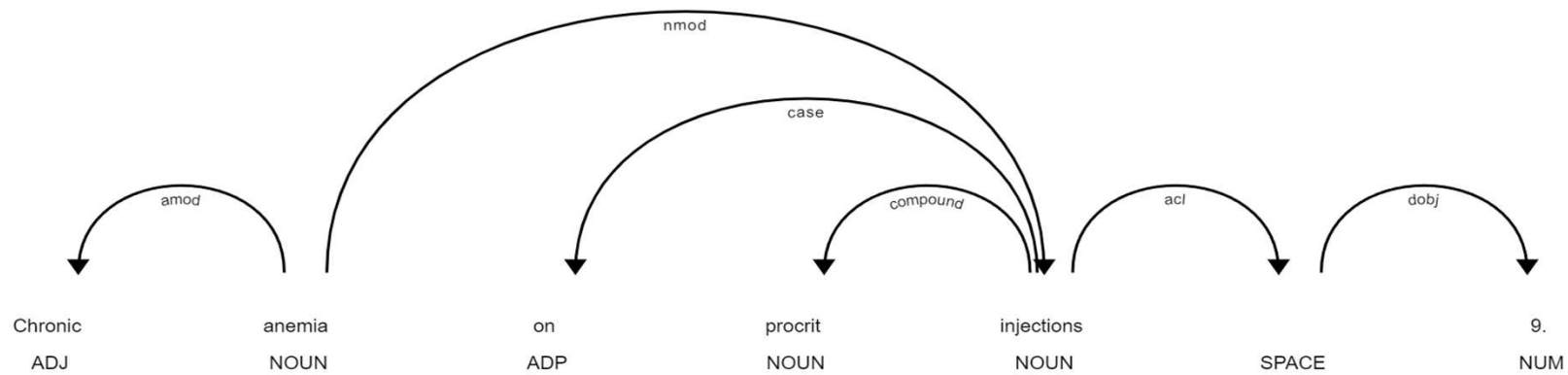
```
df =pd.read_csv("nlp_med_notes.csv")['text']
label_df=do_data_process(df,'REVIEW OF SYSTEMS')

for data in label_df:
    doc = nlp(data)
    sentence_spans = list(doc.sents)
    displacy.render(sentence_spans, style="dep", jupyter=True)
print('*****')
```

- Create a **dependency tree** using displacy from spacy model to show the relationship between words in a sentence.
- A meaningful sentence can be extracted from 'REVIEW OF SYSTEMS' sub paragraph.





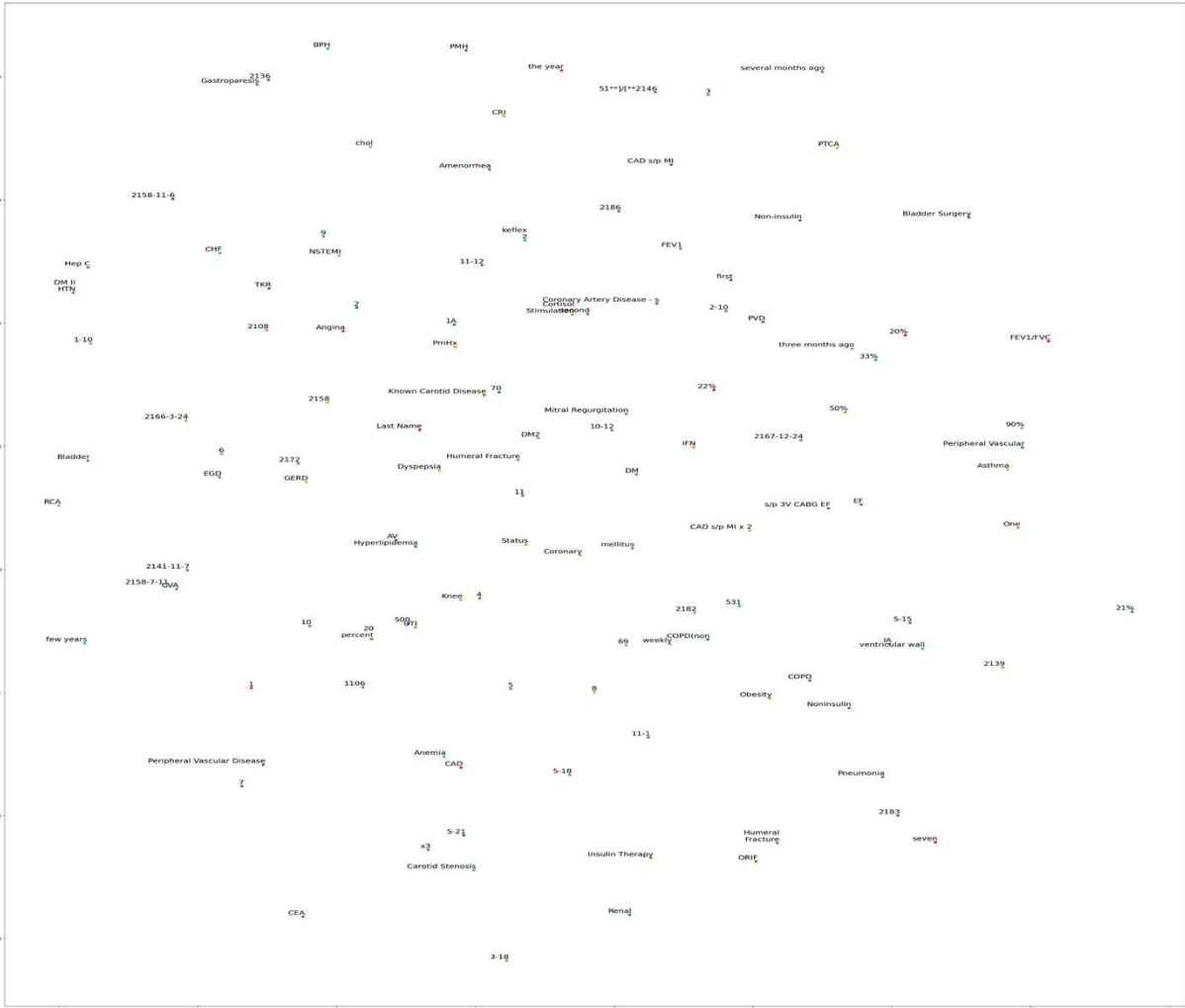


Gout  
NOUN

## Spacy TSNE visualization: For Medication in medical notes

[https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/spacy\\_tsne\\_plot\\_diag.png](https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/spacy_tsne_plot_diag.png)

```
df = pd.read_csv("nlp_med_notes.csv")['text']
label_df = do_data_process(df, 'MEDICATIONS')
```



## SciSpacy Usage:

- SciSpacy models:
  - en\_core\_sci\_md
  - en\_core\_sci\_lg
  - en\_ner\_craft\_md
  - en\_ner\_jnlpba\_md
  - en\_ner\_bionlp13cg\_md
  - en\_ner\_bc5cdr\_md
- Load the model and visualize the corpus.

```
import en_ner_craft_md
nlp = en_ner_craft_md.load()
visualize_entities(label_df)
```

## Scispacy: *en\_core\_sci\_lg*

Hepatitis C genotype 1A ENTITY complicated ENTITY by

cirrhosis ENTITY, esophageal varices ENTITY, encephalopathy ENTITY and presumptive ENTITY

SBP ENTITY, Type 2 diabetes ENTITY, Obesity ENTITY, Hypertension ENTITY, Asthma ENTITY.

Esophageal candidiasis ENTITY, Gastroparesis ENTITY, Depression ENTITY, Status

post cholecystectomy ENTITY, Status ENTITY, post ENTITY, seven, spur surgeries ENTITY.

Hypothyroidism ENTITY, Amenorrhea ENTITY, Migraines ENTITY.

\*\*\*\*\*

PmHx ENTITY :

HCV ENTITY genotype ENTITY IA refractory to IFN ENTITY x 2 ascites ENTITY

grade I esophageal varices ENTITY ( EGD ENTITY [\*\*11-1\*\*])

h/o esophageal candidiasis

s/p ccx

DM II

HTN

asthma

hypothyroid

depression

amenorrhea

migraines ENTITY

\*\*\*\*\*

1. Hypertension ENTITY.

2. History of hepatitis C ENTITY ; on pegylated interferon ENTITY and

## Spacy:

... Hepatitis C genotype 1A ENTITY complicated ENTITY by

cirrhosis, esophageal varices ENTITY, encephalopathy ENTITY and presumptive SBP ENTITY, Type 2 diabetes.

Obesity ENTITY, Hypertension ENTITY, Asthma.

Esophageal candidiasis ENTITY, Gastroparesis, Depression, Status

post cholecystectomy ENTITY, Status, post ENTITY, seven, spur surgeries ENTITY.

Hypothyroidism ENTITY, Amenorrhea, Migraines ENTITY.

... \*\*\*\*\*

... PmHx ENTITY :

HCV ENTITY genotype ENTITY IA refractory ENTITY to IFN ENTITY x 2 ascites grade I esophageal

varices ENTITY ( EGD ENTITY [\*\*11-1\*\*])

h/o esophageal candidiasis ENTITY

s/p ccx

DM II

HTN ENTITY

asthma

hypothyroid

depression

amenorrhea

migraines

... \*\*\*\*\*

... 1. Hypertension.

2. History of hepatitis C ENTITY ; on pegylated ENTITY interferon ENTITY and

Scipacy: *en\_ner\_craft\_md*

Spacy:

Hepatitis C genotype so 1A complicated by

cirrhosis, esophageal varices, encephalopathy and presumptive

SBP. Type 2 diabetes. Obesity. Hypertension. Asthma.

Esophageal candidiasis. Gastroparesis. Depression. Status

post cholecystectomy. Status post seven spur surgeries.

Hypothyroidism. Amenorrhea. Migraines.

Hepatitis C genotype 1A ENTITY complicated ENTITY by

cirrhosis, esophageal varices ENTITY, encephalopathy ENTITY and presumptive SBP ENTITY. Type 2 diabetes.

Obesity ENTITY. Hypertension ENTITY. Asthma.

Esophageal candidiasis ENTITY. Gastroparesis. Depression. Status

post cholecystectomy ENTITY. Status post ENTITY seven spur surgeries ENTITY.

Hypothyroidism ENTITY. Amenorrhea. Migraines ENTITY.

2. History of hepatitis C; on pegylated interferon CL and

acute infarct noted on MRI [\*\*5-15\*\*] (left posterior frontal region so

Congestive Heart Failure, NSTEMI, Coronary Artery Disease - s/p

multiple RCA stents, Mitral Regurgitation, Diabetes Mellitus -

on Insulin GGP Therapy, Hypercholesterolemia, Cerebrovascular

Disease - s/p CVA, Known Carotid Disease, Right Subclavian

Stenosis, Peripheral Vascular Disease, History of Humeral

Fracture, GERD, Depression, Prior Bladder Surgery

(b) Persantine CHEBI MIBI on [\*\*2183-2-4\*\*] demonstrated a

dilated left ventricle, partially reversible moderate

perfusion defect in the anterior left ventricular wall, with

moderate fixed deficits at the septal/inferior walls, global

and left ventricular hypokinesis with an ejection fraction of

21%.

- Hepatitis C cirrhosis and hepatocellular carcinoma s/p

radiofrequency ablation x 3, s/p liver transplantation [\*\*1-10\*\*]

- Recurrent Hep C after Transplant- last viral TAXON load 69 on [\*\*2158-7-11\*\*].

Scipacy: *en\_ner\_craft\_md*

- **SO** stands for Sequence Ontology
- **CL** stands for Cell Ontology
- **CHEBI** stands for Chemical Entities of Biological Interest
- **TAXON** stands for Taxonomy
- **GCP** stands for Gene Ontology Cellular Component

Scipacy: *en\_ner\_jnlpba\_md*

Hepatitis C genotype 1A **DNA** complicated by  
grade I esophageal varices ( **EGD** **CELL\_TYPE** [ **\*\*11-1\*\*** ] **DNA** )  
2. History of hepatitis C; on pegylated **interferon** **PROTEIN** and



Scipacy: *en\_ner\_bionlp13cg\_md*

1. Hypertension.

2. History of hepatitis C ORGANISM ; on pegylated interferon and ribavirin study. The patient ORGANISM quit three months ago for unknown reasons.

3. Dyspepsia/peptic ulcer CANCER disease.

4. Diabetes; per OMR CANCER but the patient ORGANISM denies.

5. Depression; no prior suicide attempts with few years of treatment.

6. Renal stones; status post renal ORGAN surgery with blood ORGANISM\_SUBSTANCE transfusions

7. History of angina.

8. Chronic back pain; status post surgery.

\*\*\*\*\*

CAD CANCER , status post CABG with an EF CELL of 20 percent.

\*\*\*\*\*

- s/p [\*\*Year (4 digits) 500\*\*] graft TISSUE from right hip to elbow

\*\*\*\*\*

AS HTN CELL elev. chol SIMPLE\_CHEMICAL .

1. Coronary artery MULTI\_TISSUE\_STRUCTURE disease, status post CABG in [\*\*2136\*\*].

2. Status post catheterization in [\*\*2141-11-7\*\*] with patent graft TISSUE .

3. Ischemic cardiomyopathy with an EF CELL of 22% with basal mid inferior, posterior IMMATERIAL\_ANATOMICAL\_ENTITY hypokinesis.

4. Bilateral carotid artery MULTI\_TISSUE\_STRUCTURE stenosis, status post left CEA GENE\_OR\_GENE\_PRODUCT .

5. Left parietal MULTI\_TISSUE\_STRUCTURE CVA.

6. Type 2 diabetes mellitus.

7. Chronic renal ORGAN insufficiency.

8. Hypercholesterolemia.

9. Hypertension.

10. Status post cholecystectomy and total abdominal hysterectomy.

11. GERD, positive H. pylori ORGANISM .

PAST MEDICAL HISTORY:

1. ESRD secondary to hypertensive nephrosclerosis s/p right upper extremity AV graft 9'[\*\*56\*\*][\*\*33\*\*] in preparation for dialysis. Graft placement was complicated by cellulitis PATHOLOGICAL\_FORMATION , for which he was

Scipacy: *en\_ner\_bc5cdr\_md*

## Conclusion:

With the help of different spacy models we are able to extract different kinds of NERs.

Hepatitis C genotype 1A complicated by

cirrhosis DISEASE , esophageal varices DISEASE , encephalopathy DISEASE and presumptive

SBP. Type 2 diabetes DISEASE . Obesity DISEASE . Hypertension DISEASE . Asthma.

Esophageal candidiasis DISEASE . Gastroparesis DISEASE . Depression DISEASE . Status

post cholecystectomy. Status post seven spur surgeries.

Hypothyroidism DISEASE . Amenorrhea DISEASE . Migraines DISEASE .

2. History of hepatitis C DISEASE ; on pegylated interferon and ribavirin CHEMICAL study.\*

PmHx: HCV genotype 1A refractory CHEMICAL to IFN CHEMICAL x 2 ascites

grade I esophageal varices DISEASE (EGD [\*\*11-1\*\*])

h/o esophageal candidiasis DISEASE

s/p ccx DM II CHEMICAL HTN

asthma

hypothyroid DISEASE

depression DISEASE

amenorrhea DISEASE

migraines DISEASE



## Word2Vec and TSNE plot usage:

Apply **fetch\_corpus()** to retrieve the list of corpus from the medical note.

```
corpus=[]
for data in label_df:
    corpus.append(fetch_corpus(data))
print(list(corpus))
```

[112] ✓ 3.6s Python

... ['diabetes mellitus', 'Hypertension', 'Chronic cervical spine disease', 'Congestive heart failure']]

```
[['cirrhosis', 'esophageal varices', 'encephalopathy', 'diabetes', 'Obesity',
'Hypertension', 'Esophageal candidiasis', 'Gastroparesis', 'Depression', 'Hypothyroidism',
'Amenorrhea', 'Migraines'], ['PmHx:\nHCV genotype IA refractory', 'IFN', 'esophageal
varices', 'esophageal candidiasis', 'ccx\ndm II\n', 'hypothyroid', 'depression',
'amenorrhea', 'migraines'], ['Hypertension', 'hepatitis C', 'pegylated interferon
and\nribavirin', 'ulcer disease', 'Diabetes', 'Depression', 'Renal stones', 'angina',
'Chronic back pain'], ['CAD'], ['Hypertension'], ['chol', 'NIDDM', 'diverticulosis', 'hiatal
hernia', 'obesity', 'appy'], ['Coronary artery disease', 'artery stenosis', 'CVA', 'diabetes
mellitus', 'Chronic renal insufficiency', 'Hypercholesterolemia', 'Hypertension', 'GERD'],
['Coronary artery disease', 'diabetes mellitus', 'Hypertension', 'Arthritis'], ['ESRD',
'hypertensive nephrosclerosis', 'cellulitis', 'keflex', 'DM', 'glyburide', 'glipizide',
'HTN', 'clonidine', 'lisinopril', 'nifedipine', 'PVD', 'CVA', 'Secondary
hyperparathyroidism', 'anemia', 'procrit', 'Prostate', 'Lupron', 'Gout'], ['ESRD',
'hypertensive nephrosclerosis', 'cellulitis', 'keflex', 'DM', 'glyburide', 'glipizide',
'HTN', 'clonidine', 'lisinopril', 'nifedipine', 'PVD', 'CVA', 'Secondary
hyperparathyroidism', 'anemia', 'procrit', 'Prostate', 'Lupron', 'Gout'], ['GERD'], ['CAD',
'CPD', 'hyperlipidemia', 'claustrophobia', 'diabetes\nmellitus type 2'], ['DM2', '^chol',
'hypothyroid', 'arthritis'], ['Coronary artery disease', 'Persantine', 'septal/inferior',
'left ventricular hypokinesis', 'Congestive heart failure', 'depressed', 'diabetes
mellitus', 'Hypertension', 'Hyperlipidemia', 'Onychomycosis', 'Anemia', 'leukocytosis',
'Chronic obstructive pulmonary disease', 'obstructive/restrictive deficit'], ['CAD', 'AICD',
'hypothyroid', 'DM', 'varicose vein removal'], ['CAD', 'atrial fibrillation', 'htn', 'GERD',
'Anemia'], ['artery disease', 'Atrial fibrillation', 'Hypertension', 'Hyperlipidemia',
'Anemia'], ['artery disease', 'Atrial fibrillation', 'Hypertension', 'Hyperlipidemia',
'Anemia'], ['Oligodendroglioma', 'oligoastrocytoma', 'infertility', 'temozolomide',
'seizures', 'temozolomide', 'dexamethasone', 'sepsis', 'encephalopathy', 'tumor',
'weakness', 'steroid myopathy', 'Hyperglycemia', 'steroid'], ['UTIs', 'NIDDM',
'Hypercholesterolemia', 'Autoimmune Hepatitis'], ['cirrhosis', 'HCC', 'cirrhosis',
'Ascites', 'encephalopathy', 'HD', 'MMF'], ['cirrhosis and hepatocellular carcinoma',
'liver\nfailure', 'ascites', 'encephalopathy', 'Type II DM\n- Adrenal Insufficiency:
**2158-11-6*'], ['Cortisol', 'Urolithiasis'], ['Hypertension'], ['cirrhosis and
hepatocellular carcinoma'], ['Diabetes', 'Dyslipidemia', 'Hypertension', 'prostatic
hypertrophy', 'Arthritis', 'Gout', 'Bladder stone'], ['dilated cardiomyopathy', 'mitral
regurgitation', 'diabetes', 'steroids', 'Pneumonia', 'ceftriaxone', 'azithromycin'],
['Hypertension', 'DM II', 'CAD', 'steroids', 'duodenal ulcer', 'CHF', 'dementia']]
```

## Code Snippet:

The resultant presents the similarity logits for the word 'encephalopathy'

```
model1 = Word2Vec(corpus, min_count=1)
model1.wv['encephalopathy']
```

✓ 0.0s

```
array([-8.7531786e-03,  2.1741530e-03, -8.6094369e-04, -9.3106795e-03,
       -9.4064260e-03, -1.4538610e-03,  4.4581434e-03,  3.7536507e-03,
       -6.5508662e-03, -6.8758638e-03, -5.0241956e-03, -2.3389754e-03,
       -7.2221956e-03, -9.5775630e-03, -2.7493779e-03, -8.3579253e-03,
       -6.0137236e-03, -5.7307631e-03, -2.3647691e-03, -1.7745970e-03,
       -8.9362776e-03, -6.9540367e-04,  8.1498129e-03,  7.6987552e-03,
       -7.2270953e-03, -3.6619261e-03,  3.0702241e-03, -9.5547633e-03,
        1.4801961e-03,  6.5093059e-03,  5.7971054e-03, -8.7778289e-03,
       -4.5126704e-03, -8.1743700e-03,  3.6444104e-05,  9.3236137e-03,
        5.9737498e-03,  5.0418158e-03,  5.0477851e-03, -3.3005176e-03,
        9.5378207e-03, -7.3622023e-03, -7.3122410e-03, -2.2796686e-03,
       -7.5115956e-04, -3.1877523e-03, -6.4017362e-04,  7.4983523e-03,
       -6.7837693e-04, -1.5929459e-03,  2.7603914e-03, -8.3850855e-03,
        7.8556603e-03,  8.5417535e-03, -9.6132429e-03,  2.4651806e-03,
        9.9031590e-03, -7.6433863e-03, -6.9885631e-03, -7.6803914e-03,
        8.3996654e-03, -6.9426361e-04,  9.1576520e-03, -8.1540635e-03,
        3.7199876e-03,  2.6663735e-03,  7.5992203e-04,  2.3442844e-03,
       -7.5090886e-03, -9.2971604e-03,  2.3168572e-03,  6.1675226e-03,
        8.0000097e-03,  5.6976336e-03, -7.6059706e-04,  8.2836589e-03,
       -9.3513643e-03,  3.3959236e-03,  2.5762038e-04,  3.8506044e-03,
        7.3216450e-03, -6.7115389e-03,  5.5358177e-03, -9.4783595e-03,
       -8.4873615e-04, -8.6890254e-03, -5.0572841e-03,  9.3041677e-03,
       -1.8036201e-03,  2.8908700e-03,  9.0945661e-03,  8.9400755e-03,
       -8.2035558e-03, -3.0187166e-03,  9.9292845e-03,  5.0835693e-03,
       -1.5810047e-03, -8.7010888e-03,  2.9685348e-03, -6.6635790e-03],
      dtype=float32)
```

The code snippet shows the similar words that the model contains within it.

```
model1.wv.similar_by_word('encephalopathy')
```

✓ 0.0s

```
[('diverticulosis', 0.30835863947868347),  
 ('Renal stones', 0.2804599106311798),  
 ('esophageal candidiasis', 0.2360624074935913),  
 ('Hyperglycemia', 0.20393291115760803),  
 ('Fracture', 0.2022656798362732),  
 ('Persantine', 0.19107292592525482),  
 ('Gastroparesis', 0.17959770560264587),  
 ('Hypercholesterolemia', 0.168697327375412),  
 ('Secondary hyperparathyroidism', 0.16342854499816895),  
 ('Peripheral Vascular Disease', 0.14657379686832428)]
```

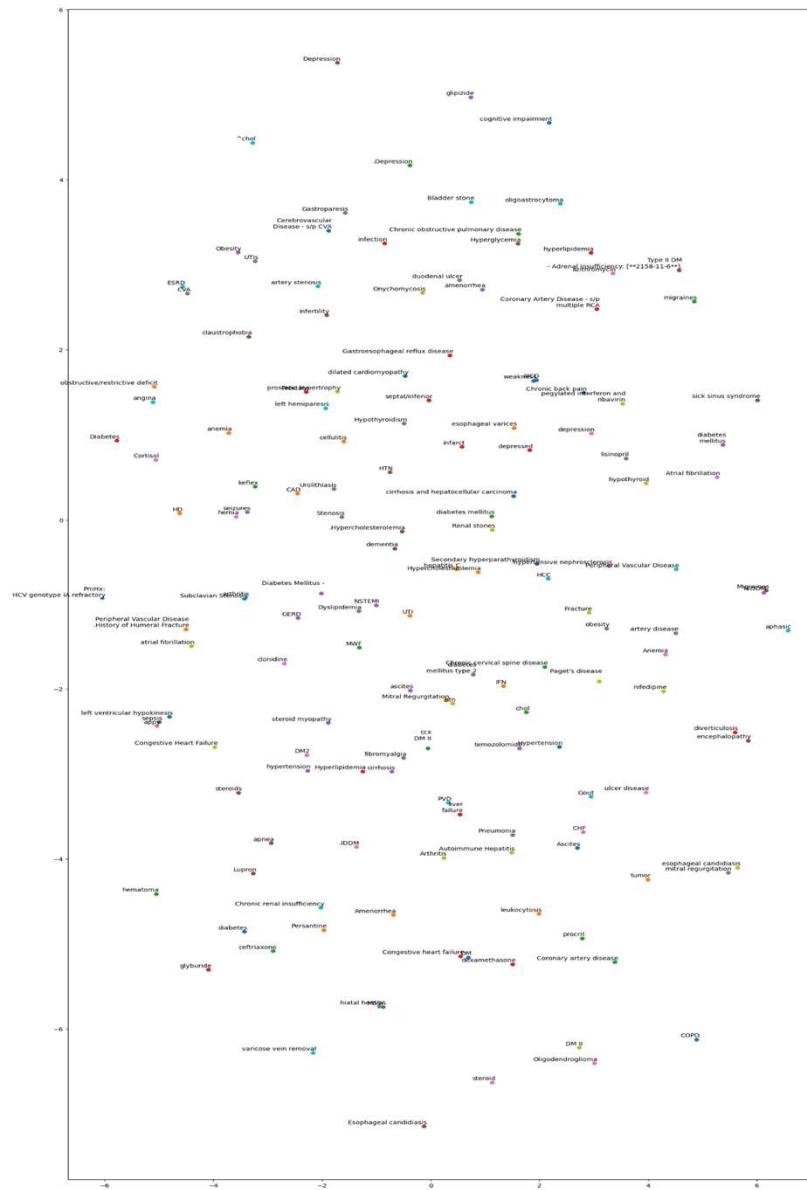
### TSNE plot:

```
vocabs = model1.wv.key_to_index.keys()
new_v = list(vocabs)
tsne_plot(model1, new_v)
```

The corpus are plotted in tsne word embedding as per their logits.

### ***SciSpacy TSNE plot:***

[https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/tsne\\_plot\\_diag.png](https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/tsne_plot_diag.png)





## Clinical BERT:

- Import transformers and classification package.

```
from transformers import AutoTokenizer,  
AutoModelForTokenClassification  
from transformers import pipeline
```

- Initialize model.

```
tokenizer = AutoTokenizer.from_pretrained("emilyalsentzer/Bio_ClinicalBERT")  
model = AutoModelForTokenClassification.from_pretrained("emilyalsentzer/Bio_ClinicalBERT")  
nlp = pipeline("ner", model=model, tokenizer=tokenizer)
```

- Load dataset and examine 'PHYSICAL EXAMINATION' sub passage.

- The output displays label and score for each word in a sentence.
- The label is categorized by ClinicalBERT into LABEL\_0 and LABEL\_1
- The score is given to each word in that label category.

```
result=[]
for data in label_df:
    value=nlp(data)
    print(value)
```

✓ 0.4s

```
[{'entity': 'LABEL_0', 'score': 0.6019393, 'index': 1, 'word': 'the', 'start': 0, 'end': 3}, {
[{'entity': 'LABEL_0', 'score': 0.5833172, 'index': 1, 'word': 'the', 'start': 0, 'end': 3}, {
[{'entity': 'LABEL_1', 'score': 0.524752, 'index': 1, 'word': 'vital', 'start': 0, 'end': 5},
[{'entity': 'LABEL_1', 'score': 0.59145355, 'index': 1, 'word': 'vs', 'start': 0, 'end': 2}, {
[{'entity': 'LABEL_1', 'score': 0.60578763, 'index': 1, 'word': 'height', 'start': 0, 'end': 6
[{'entity': 'LABEL_0', 'score': 0.5948016, 'index': 1, 'word': '-', 'start': 0, 'end': 1}, {'e
[{'entity': 'LABEL_1', 'score': 0.52503663, 'index': 1, 'word': 'vital', 'start': 0, 'end': 5}
[{'entity': 'LABEL_1', 'score': 0.5799533, 'index': 1, 'word': 'vital', 'start': 0, 'end': 5},
```

- For tokenization use import nltk package
- nltk.tokenize is used to split each text in list into individual words

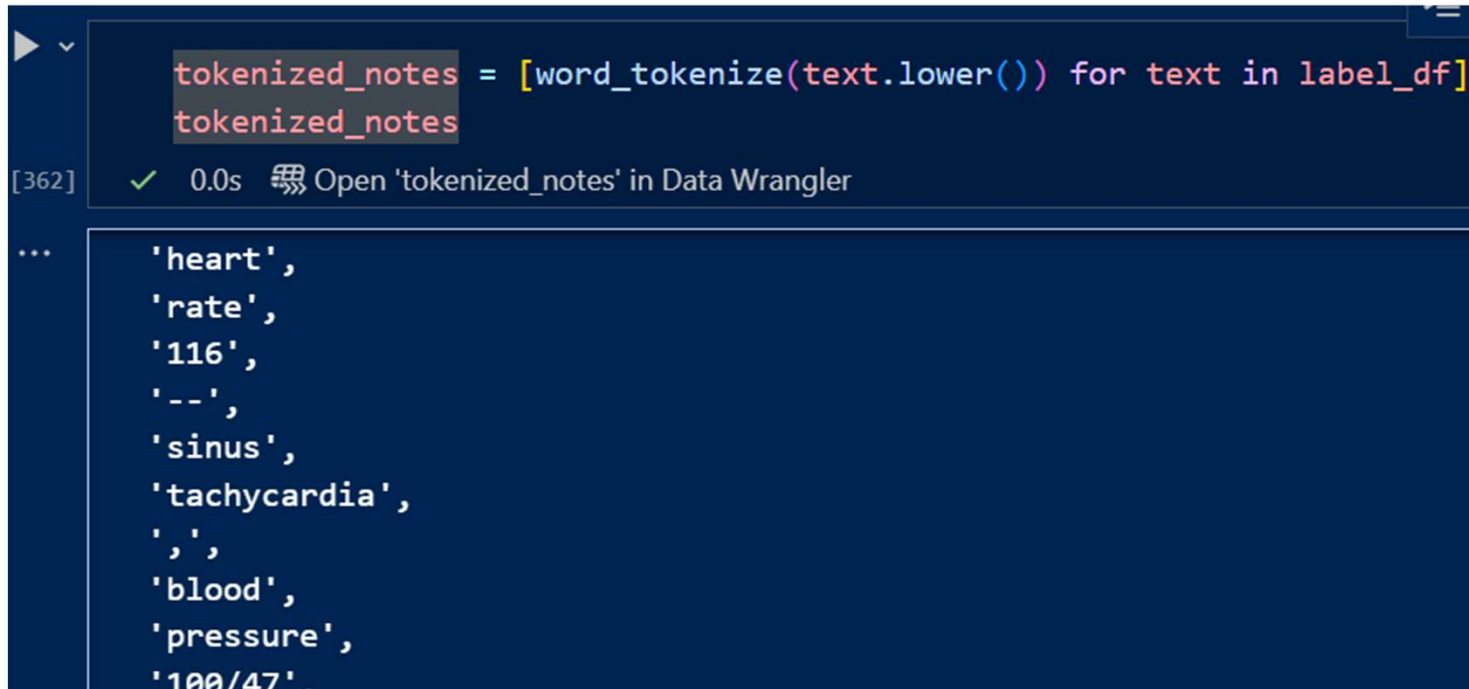
```
from nltk.tokenize import word_tokenize
import nltk
nltk.download('punkt')
```



- ClinicalBert can be visualized on 'PHYSICAL EXAM' label extracted subpassage from medical notes.

```
df = pd.read_csv("nlp_med_notes.csv")['text']  
label_df = do_data_process(df, 'Physical Exam')
```

- Tokenized notes presents the tokens extracted from the 'PHYSICAL EXAM' label\_df



The screenshot shows a Jupyter Notebook interface. The top part displays a code cell with the following Python code:

```
tokenized_notes = [word_tokenize(text.lower()) for text in label_df]  
tokenized_notes
```

Below the code cell, the output is shown as a list of tokens for a specific medical note. The output is displayed as a list of strings, with some tokens truncated by ellipses at the end of the line.

```
...  
'heart',  
'rate',  
'116',  
'--',  
'sinus',  
'tachycardia',  
',',  
'blood',  
'pressure',  
'100/47',
```

- Word Vector is created by substituting the tokenized\_notes.

```
model = Word2Vec(sentences=tokenized_notes, vector_size=100, window=5, min_count=2,  
workers=4)
```

**ClinicalBERT TSNE visualization**: For Medication in medical notes

[https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/clinical\\_bert\\_tsne\\_plot.png](https://github.com/dalwari/mimic-iii-clinical-database-demo-1.4/blob/main/clinical_bert_tsne_plot.png)

