

# CSE4022 – Natural Language Processing Final Report

## Biomedical Text Mining (Acronym and Text Extraction) Doctor Helper

Parth Maheshwari – 19BCT0221  
Pramit Gupta – 19BCT0204  
Rohit Ravichandran – 19BCT0014

## Abstract

Doctor-helper is a two way information extraction/biomedical text-mining project. It uses Python as a base and many NLP related libraries such as spaCy for extracting information. It uses the standard 6-step process to process raw unstructured data into formattable structured data. First it takes input from the patient (user1-patient), processes and converts it into a tabular format with some mandatory fields such as patient name, age and contact number, and other fields like address, symptoms, body-temperature, etc which is achieved by a modified parts of speech tagging. It will also contain a field which displays the possible ailments which the patient could be suffering from, based on probability calculated using symptoms observed, weather of the address, time of year, last known cases/outbreaks. The doctor can use this information and it will be stored in a patient-database. In the second stage, the doctor will write a report which will be taken as input (user2-doctor), processed and converted it into a tabular format which is easier to read and be understood by the patient. It will contain fields like ailments(confirmed and suspected), tests to be taken(blood test, etc.), drugs prescribed, dosage, side effects, alternative drugs, etc. It will also give the patient an option to add a reminder for each dosage using the Google Calender API.

## Introduction

In these trying times, the importance of doctors worldwide has risen and so have their responsibilities. COVID-19, being highly contagious, has significantly brought down the accessibility of medical help to patients suffering from ailments other than COVID-19. Whether it be due to medical professionals being engaged by COVID-19 patients or due to the person being afraid to seek help from a clinic/hospital from fear of contracting the infamous corona virus, this has directly affected the patients who have contracted other ailments.

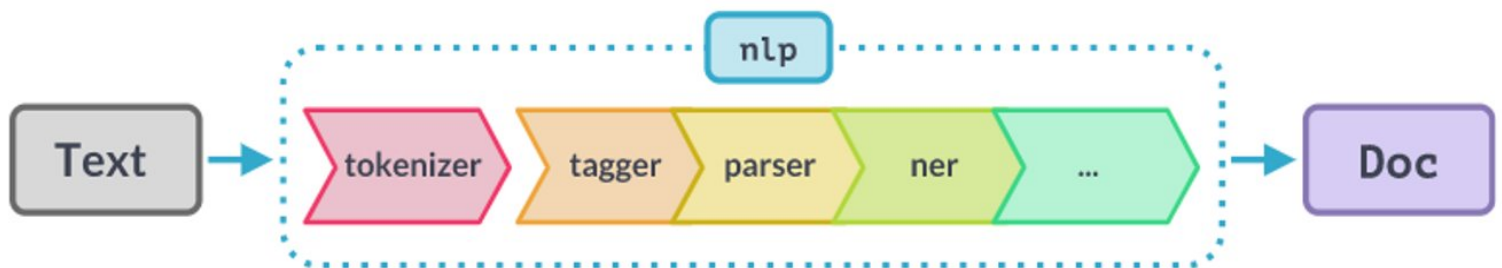
Doctor-helper aims to ease the load on medical professionals by making their task much easier. It takes input from the patient and extracts the patient details, various symptoms present, etc. and detects the possible ailments that the patient may be suffering from (which is then confirmed by the doctor). Then it takes the doctor's report as input and displays the output in a tabular format, containing the drugs prescribed, dosage, precautions, side-effects, alternative drugs, etc. It will also give an option to set a reminder for each dose using the Google Calendar.

This application will help streamline the process of consulting the doctor and save their time on unnecessary paperwork by automating most of it. It will aid the doctor's judgement and also help with online consultation from medical professionals without the patient risking contact with individuals affected from COVID-19 and following lockdown mandates.

## Problem Statement

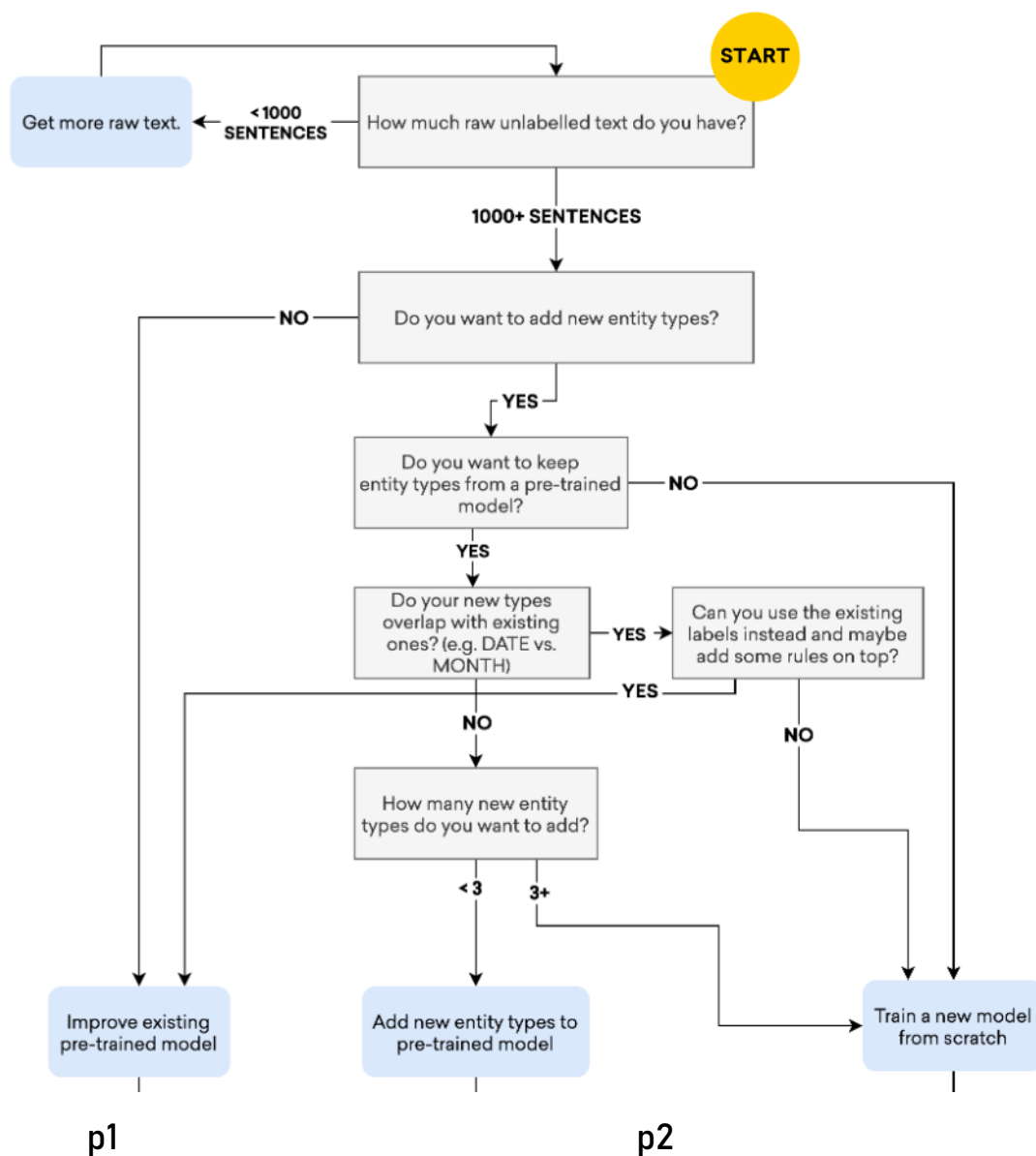
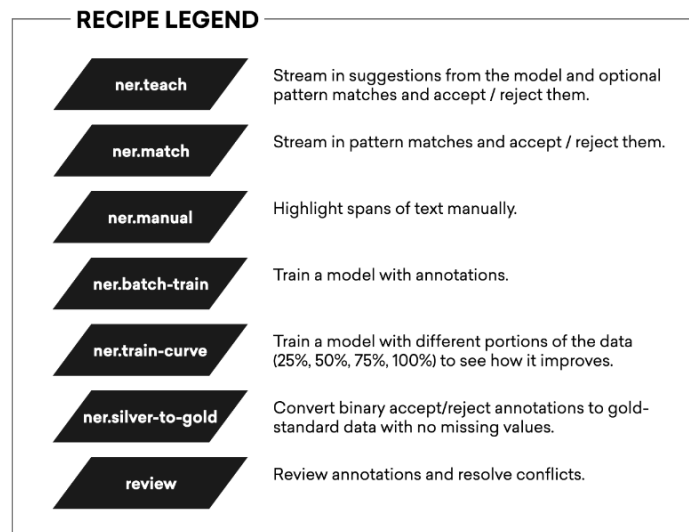
- We are expecting to implement basic concepts of Text Mining and NLP in this review.
- We are going to play with the small CORPUS i.e. data-set we have right now.
- We have made sure that the data set is complicated enough, i.e. we have taken random data from different journals.
- In this review we are going to label the DRUG and ILLNESS NAMES, that are present in the data-set.
- We are going to use MACHINE LEARNING, NLP and TEXT MINING to achieve so.

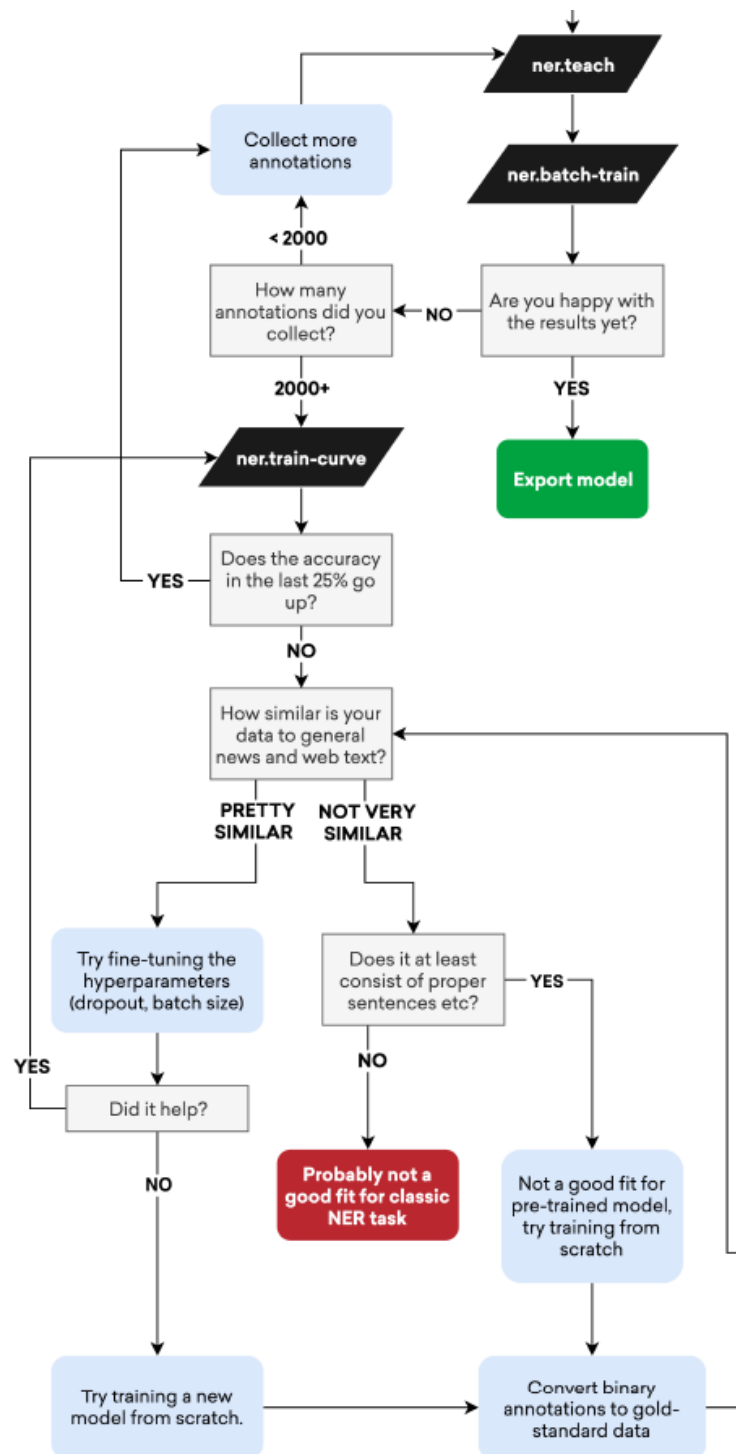
## 🌀 Architecture Diagram



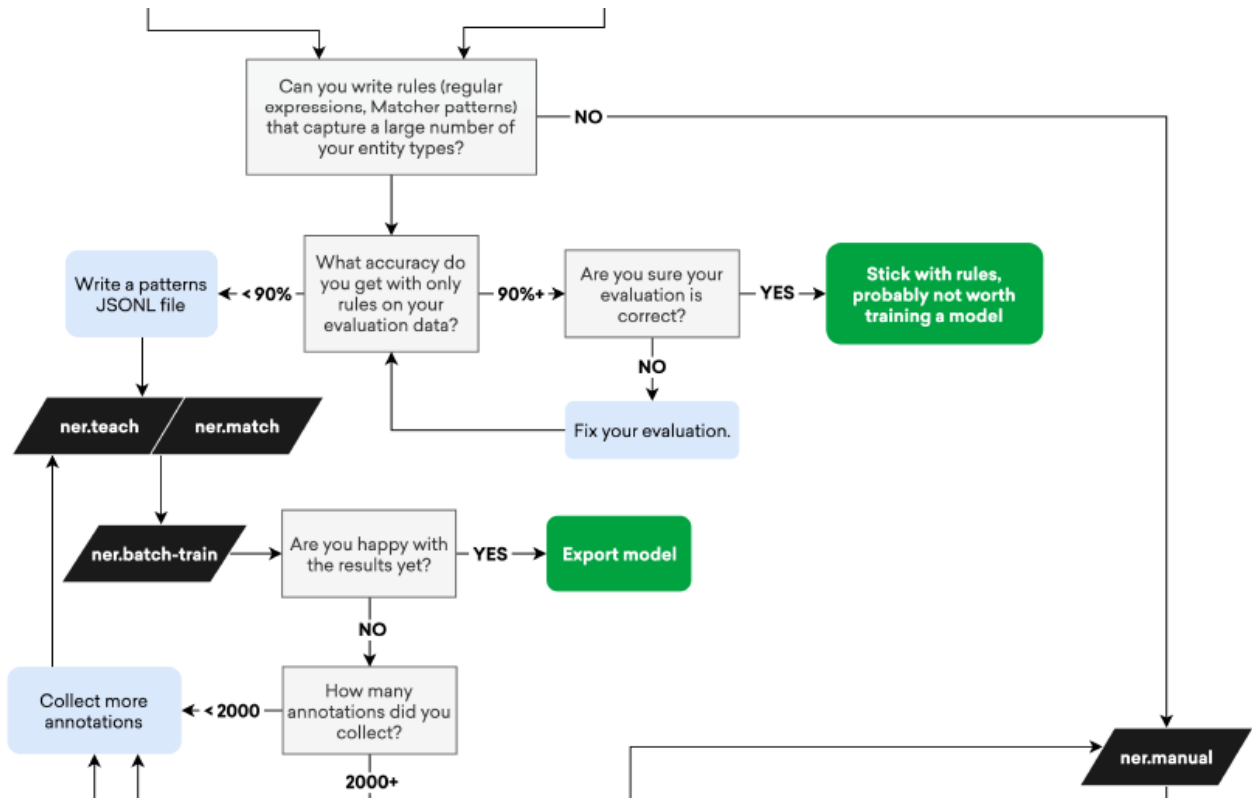
NAME	COMPONENT	CREATES	DESCRIPTION
<b>tokenizer</b>	<a href="#">Tokenizer</a> ≡	<code>Doc</code>	Segment text into tokens.
<b>PROCESSING PIPELINE</b>			
<b>tagger</b>	<a href="#">Tagger</a> ≡	<code>Token.tag</code>	Assign part-of-speech tags.
<b>parser</b>	<a href="#">DependencyParser</a> ≡	<code>Token.head</code> , <code>Token.dep</code> , <code>Doc.sents</code> , <code>Doc.noun_chunks</code>	Assign dependency labels.
<b>ner</b>	<a href="#">EntityRecognizer</a> ≡	<code>Doc.ents</code> , <code>Token.ent_iob</code> , <code>Token.ent_type</code>	Detect and label named entities.
<b>lemmatizer</b>	<a href="#">Lemmatizer</a> ≡	<code>Token.lemma</code>	Assign base forms.
<b>textcat</b>	<a href="#">TextCategorizer</a> ≡	<code>Doc.cats</code>	Assign document labels.
<b>custom</b>	<a href="#">custom components</a>	<code>Doc._.xxx</code> , <code>Token._.xxx</code> , <code>Span._.xxx</code>	Assign custom attributes, methods or properties.

# 🌀 Flow Diagram



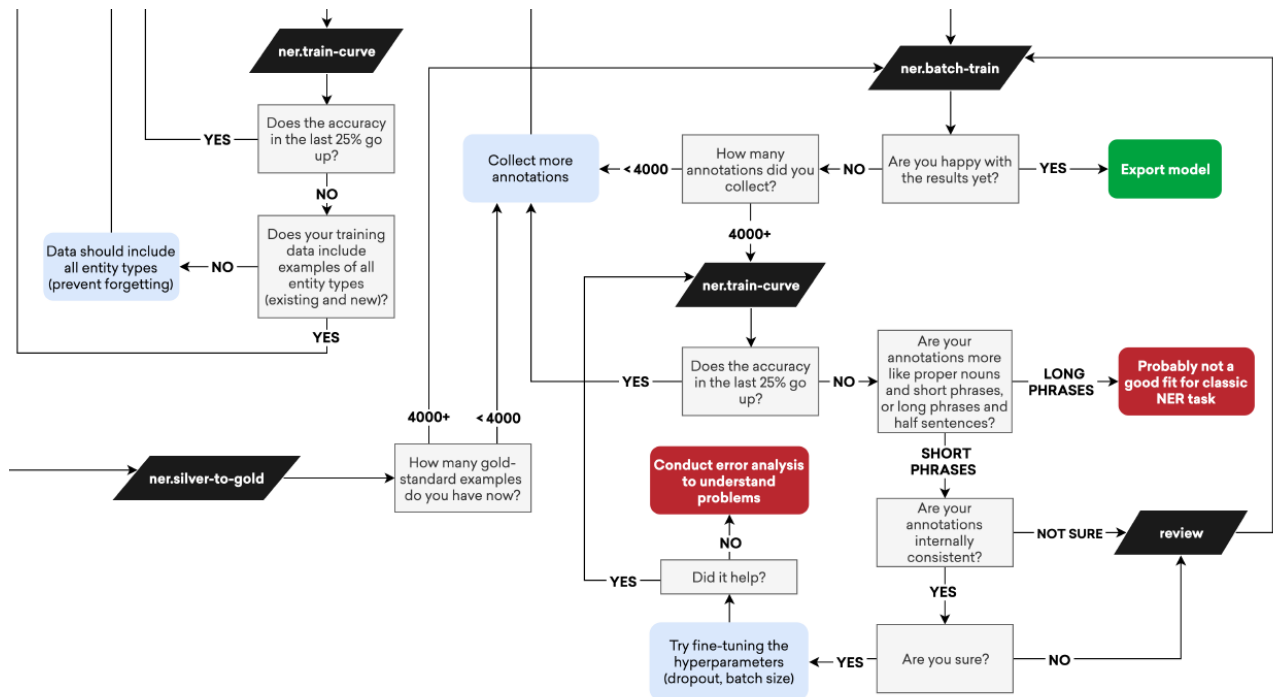


## ◆ P2





# ◆ P3



## Pseudocode

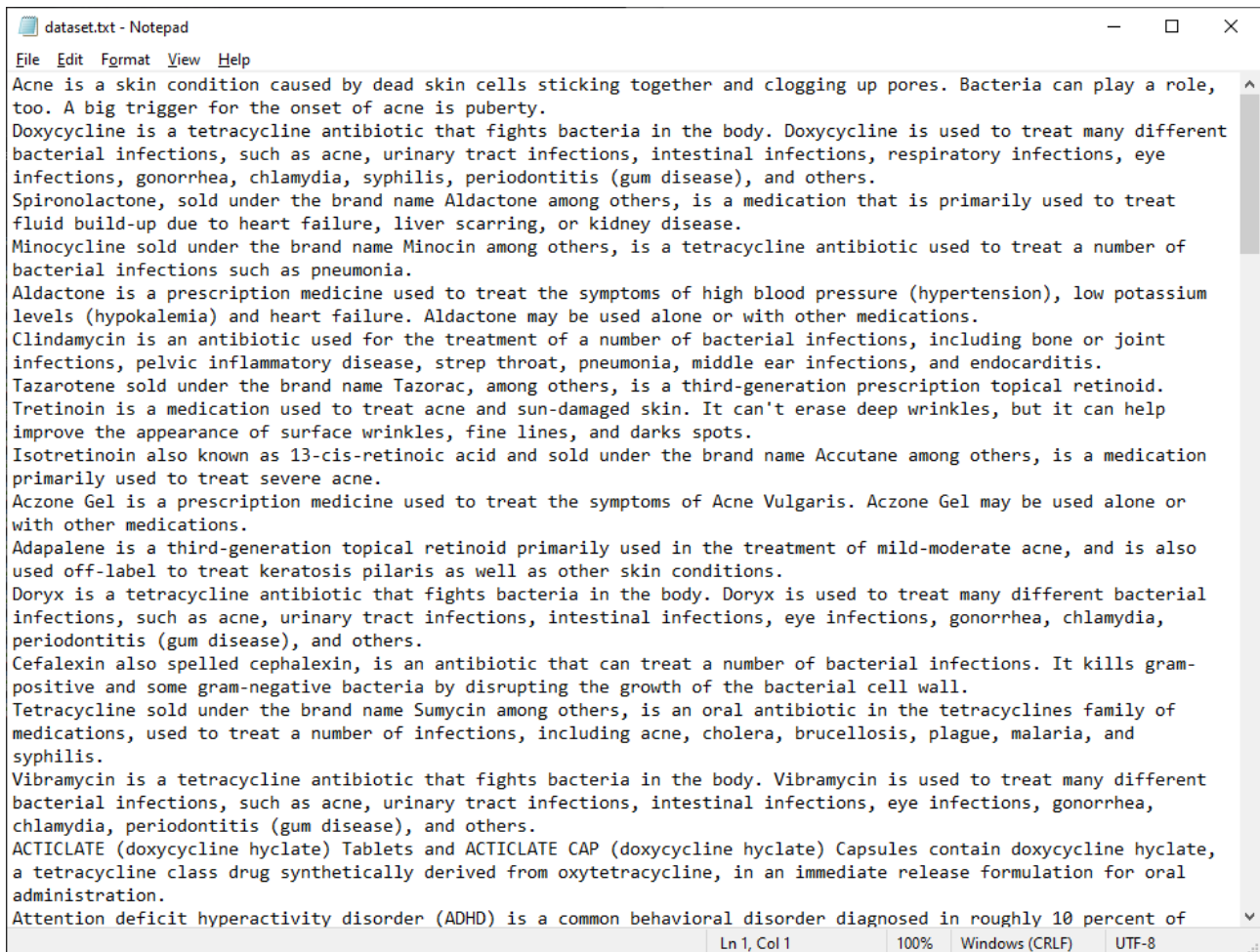
>>

```
import required packages
import data set and convert to data frame
create object for spaCy using default english corpus
pass each sentence in data frame through object for processing
using nlp(tokenizer→tagger→parser→ner→postprocessing)
import sample drug corpus and convert to data frame
apply spaCy ner pipeline with new corpus to processed data
using loop, identify index position of each drug mentioned and tag
train data using loop and use displaCy to display result
do the same for illnesses using illness.txt corpus
```

# Experiment

## ♦ Data set

- taken from various medical journals and medical reports
- contains patient descriptions, medical prescriptions, blogs about medical health, etc.



The screenshot shows a Notepad window titled "dataset.txt - Notepad". The text inside is a collection of medical descriptions and drug information, including:

- Acne is a skin condition caused by dead skin cells sticking together and clogging up pores. Bacteria can play a role, too. A big trigger for the onset of acne is puberty.
- Doxycycline is a tetracycline antibiotic that fights bacteria in the body. Doxycycline is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, respiratory infections, eye infections, gonorrhea, chlamydia, syphilis, periodontitis (gum disease), and others.
- Spiroonolactone, sold under the brand name Aldactone among others, is a medication that is primarily used to treat fluid build-up due to heart failure, liver scarring, or kidney disease.
- Minocycline sold under the brand name Minocin among others, is a tetracycline antibiotic used to treat a number of bacterial infections such as pneumonia.
- Aldactone is a prescription medicine used to treat the symptoms of high blood pressure (hypertension), low potassium levels (hypokalemia) and heart failure. Aldactone may be used alone or with other medications.
- Clindamycin is an antibiotic used for the treatment of a number of bacterial infections, including bone or joint infections, pelvic inflammatory disease, strep throat, pneumonia, middle ear infections, and endocarditis.
- Tazarotene sold under the brand name Tazorac, among others, is a third-generation prescription topical retinoid.
- Tretinoin is a medication used to treat acne and sun-damaged skin. It can't erase deep wrinkles, but it can help improve the appearance of surface wrinkles, fine lines, and dark spots.
- Isotretinoin also known as 13-cis-retinoic acid and sold under the brand name Accutane among others, is a medication primarily used to treat severe acne.
- Aczone Gel is a prescription medicine used to treat the symptoms of Acne Vulgaris. Aczone Gel may be used alone or with other medications.
- Adapalene is a third-generation topical retinoid primarily used in the treatment of mild-moderate acne, and is also used off-label to treat keratosis pilaris as well as other skin conditions.
- Doryx is a tetracycline antibiotic that fights bacteria in the body. Doryx is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, eye infections, gonorrhea, chlamydia, periodontitis (gum disease), and others.
- Cefalexin also spelled cephalixin, is an antibiotic that can treat a number of bacterial infections. It kills gram-positive and some gram-negative bacteria by disrupting the growth of the bacterial cell wall.
- Tetracycline sold under the brand name Sumycin among others, is an oral antibiotic in the tetracyclines family of medications, used to treat a number of infections, including acne, cholera, brucellosis, plague, malaria, and syphilis.
- Vibramycin is a tetracycline antibiotic that fights bacteria in the body. Vibramycin is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, eye infections, gonorrhea, chlamydia, periodontitis (gum disease), and others.
- ACTICLATE (doxycycline hyclate) Tablets and ACTICLATE CAP (doxycycline hyclate) Capsules contain doxycycline hyclate, a tetracycline class drug synthetically derived from oxytetracycline, in an immediate release formulation for oral administration.
- Attention deficit hyperactivity disorder (ADHD) is a common behavioral disorder diagnosed in roughly 10 percent of

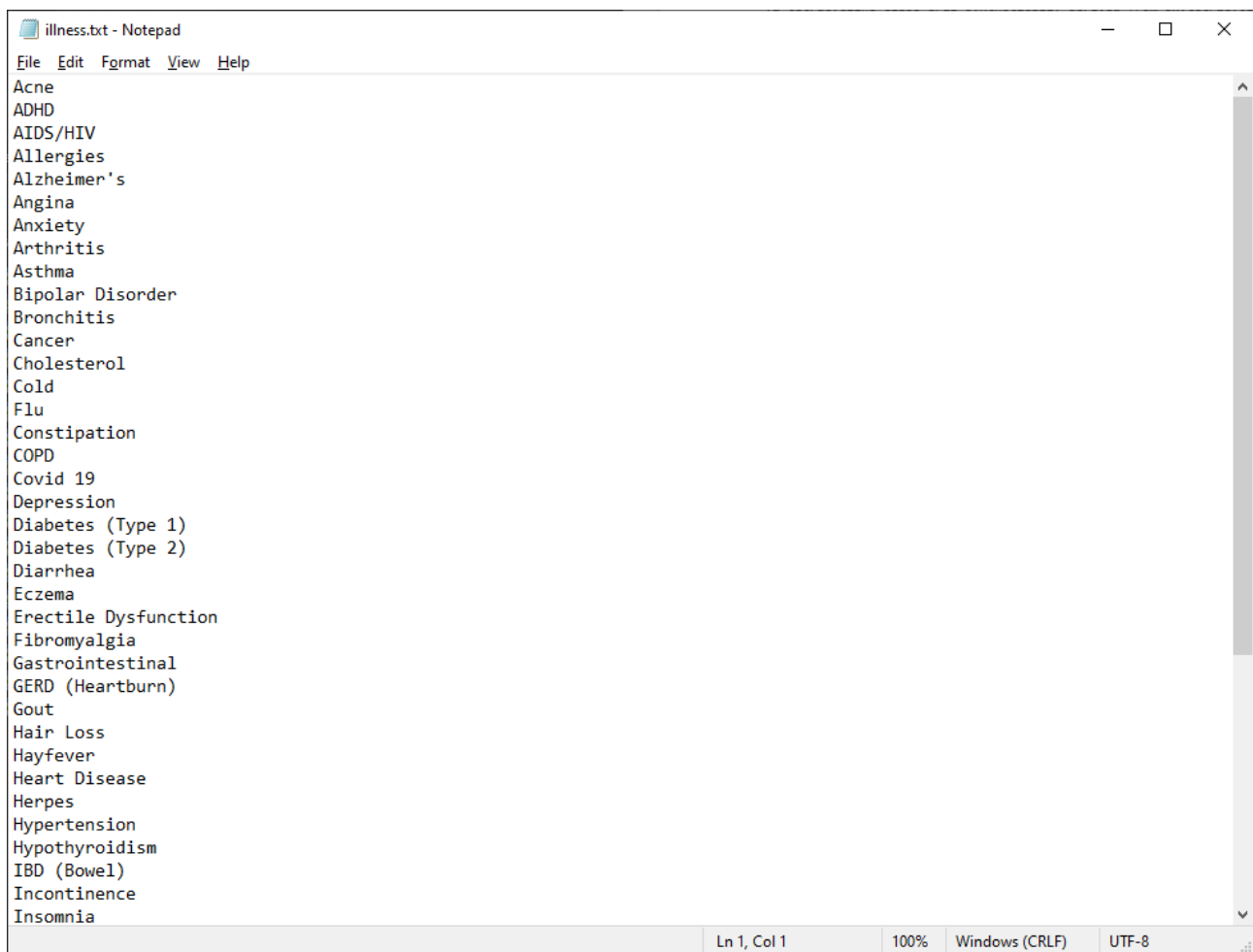
The status bar at the bottom indicates "Ln 1, Col 1", "100%", "Windows (CRLF)", and "UTF-8".

## ♦ Corpus

- corpus for ner tagging
- contains list of 20000+ pharmaceutical drugs

```
drugs.txt - Notepad
File Edit Format View Help
doxycycline
spironolactone
minocycline
Aldactone
clindamycin
Tazorac
tazarotene
tretinoin
isotretinoin
Bactrim
Epiduo
Aczone
Differin
benzoyl peroxide
adapalene
Doryx
Septra
Solodyn
cephalexin
tetracycline
Vibramycin
Acticlate
Vyvanse
Adderall
Concerta
amantadine
methylphenidate
Strattera
Dexedrine
Ritalin
clonidine
Intuniv
bupropion
guanfacine
Focalin
Desoxyn
lisdexamfetamine
Ln 1, Col 1 100% Windows (CRLF) UTF-8
```

- contains list of illness



The image shows a Notepad window with the title "illness.txt - Notepad". The menu bar includes "File", "Edit", "Format", "View", and "Help". The text area contains a list of medical conditions, one per line. The status bar at the bottom indicates "Ln 1, Col 1", "100%", "Windows (CRLF)", and "UTF-8".

```
illness.txt - Notepad
File Edit Format View Help
Acne
ADHD
AIDS/HIV
Allergies
Alzheimer's
Angina
Anxiety
Arthritis
Asthma
Bipolar Disorder
Bronchitis
Cancer
Cholesterol
Cold
Flu
Constipation
COPD
Covid 19
Depression
Diabetes (Type 1)
Diabetes (Type 2)
Diarrhea
Eczema
Erectile Dysfunction
Fibromyalgia
Gastrointestinal
GERD (Heartburn)
Gout
Hair Loss
Hayfever
Heart Disease
Herpes
Hypertension
Hypothyroidism
IBD (Bowel)
Incontinence
Insomnia
Ln 1, Col 1 100% Windows (CRLF) UTF-8
```

## Code:

```
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

import re
import random

df = pd.read_csv('sample_data/dataset.txt', sep='\t',
names=["info"])

df.head()

import spacy
from spacy import displacy
nlp = spacy.load('en_core_web_sm')

for sentence in df['info']:
    print("Sentence is: ", sentence)
    sentence_doc = nlp(sentence)
    displacy.render(sentence_doc, style='ent', jupyter=True)

df_drug = pd.read_csv('sample_data/drugs.txt', sep='\t',
names=["drug"])

df_drug.head()

df_illness = pd.read_csv('sample_data/illness.txt',
sep='\t', names=["illness"])

df_illness.head()

import spacy
from spacy.util import minibatch, compounding
```

```

nlp0 = spacy.load('en_core_web_sm')

nlp0.pipe_names

ner0 = nlp0.get_pipe('ner')

all_drugs = df_drug['drug'].unique().tolist()
all_drugs = [x.lower() for x in all_drugs]

all_drugs

all_illness = df_illness['illness'].unique().tolist()
all_illness = [x.lower() for x in all_illness]

all_illness

df['info']

def process_review(review):
    processed_token = []
    for token in review.split():
        token = ''.join(e.lower() for e in token if
e.isalnum())
        processed_token.append(token)
    return ' '.join(processed_token)

# Training examples in the required format
count = 0
TRAIN_DATA = []
for _, item in df.iterrows():
    ent_dict = {}
    if count < 1000:
        review = process_review(item['info'])
        visited_item = []
        entities = []
        for token in review.split():
            if token in all_drugs:
                for i in re.finditer(token, review):
                    if token not in visited_item:
                        entity = (i.span()[0], i.span()[1],
'DRUG')

```

```

        visited_item.append(token)
        entities.append(entity)
    elif token in all_illness:
        for i in re.finditer(token, review):
            if token not in visited_item:
                entity = (i.span()[0], i.span()[1],
'ILLNESS')
                visited_item.append(token)
                entities.append(entity)
    if len(entities) > 0:
        ent_dict['entities'] = entities
        train_item = (review, ent_dict)
        TRAIN_DATA.append(train_item)
        count+=1

```

TRAIN\_DATA

```

# Add the new label to ner
ner0.add_label(LABEL1)
ner0.add_label(LABEL2)

```

```

# Resume training
optimizer = nlp0.resume_training()
move_names = list(ner0.move_names)

```

```

# List of pipes you want to train
pipe_exceptions = ["ner", "trf_wordpiecer",
"trf_tok2vec"]

```

```

# List of pipes which should remain unaffected in
training
other_pipes = [pipe for pipe in nlp0.pipe_names if pipe
not in pipe_exceptions]

```

```

n_iter = 100
def train_ner(training_data):
    TRAIN_DATA = training_data
    nlp1 = spacy.blank("en")
    print("Created blank 'en' model")

```

```

    if "ner" not in nlp1.pipe_names:

```



```

    print("ner is created")
    ner = nlp1.create_pipe("ner")
    nlp1.add_pipe(ner, last=True)

else:
    ner = nlp1.get_pipe("ner")

for _, annotations in TRAIN_DATA:
    for ent in annotations.get("entities"):
        ner.add_label(ent[2])
        # print(ner)
nlp1.begin_training()
for itn in range(n_iter):
    random.shuffle(TRAIN_DATA)
    losses = {}
    batches = minibatch(TRAIN_DATA,
size=compounding(4.0, 32.0, 1.001))
    for batch in batches:
        texts, annotations = zip(*batch)
        nlp1.update(
            texts,
            annotations,
            drop=0.5,
            losses=losses,
        )
        print("Losses", losses)
return nlp1

```

```
nlp3 = train_ner(TRAIN_DATA)
```

```
from spacy import displacy
```

```

sentence="Paracetamol, also known as acetaminophen, is a
medication used to treat fever and mild to moderate
pain."
print("Sentence is: ", sentence)
sentence_doc = nlp3(sentence)
colors = {"DRUG": "coral", "ILLNESS": "GREEN"}
options = {"ents": ["DRUG", "ILLNESS"], "colors": colors}

```

```
displacy.render(sentence_doc, style='ent', jupyter=True, options=options)
```

```
for sentence in df['info']:
    print("Sentence is: ", sentence)
    sentence_doc = nlp3(sentence)
    colors = {"DRUG": "coral", "ILLNESS": "GREEN"}
    options = {"ents": ["DRUG", "ILLNESS"], "colors":
colors}
```

```
displacy.render(sentence_doc, style='ent', jupyter=True, options=options)
```

```
def find_terms(text):
    terms = []
    review = process_review(text)
    for token in review.split():
        if token in all_drugs:
            terms.append(token)
    return terms
```

```
# apply function
df['drug_terms'] = df['info'].apply(find_terms)
```

```
def find_terms_ill(text):
    terms = []
    review = process_review(text)
    for token in review.split():
        if token in all_illness:
            terms.append(token)
    return terms
```

```
# apply function
df['illness_terms'] = df['info'].apply(find_terms_ill)
```

```
df.head()
```

```
sensitive = pd.DataFrame(columns=['Content', 'Drug_Term'])
unsensitive = pd.DataFrame(columns=['Content'])
row_list1 = []
row_list2 = []
```

```

for i in range(len(df)):
    if len(df.loc[i,'drug_terms'])!=0:
        for k in df.loc[i,'drug_terms']:

row_list1.append({'Content':df.loc[i,'info'],'Drug_Term':
k})
    else:
        row_list2.append({'Content':df.loc[i,'info']})
sensitive = pd.DataFrame(row_list1)
unsensitive = pd.DataFrame(row_list2)

sensitive.head()

unsensitive.head()

```

```

drug_suffix = {"azole":"antifungal (except
metronidazole)",
"caine":"anesthetic",
"cillin":"antibiotic(penicillins)",
"mycin":"antibiotic",
"micin":"antibiotic",
"cycline":"antibiotic",
"oxacin":"antibiotic",
"ceph":"antibiotic(cephalosporins)",
"cef":"antibiotic (cephalosporins)",
"dine":"h2 blockers (anti-ulcers)",
"done":"opiod analgesics",
"ide":"oral hypoglycemics",
"lam":"anti-anxiety",
"pam":"anti-anxiety",
"mide":"diuretics",
"zide":"diuretics",
"nium":"neuromuscular blocking agents",
"olol":"beta blockers",
"tidine":"h2 antagonist",
"tropin":"pituitary hormone",
"zosin":"alpha blocker",
"ase":"thrombolytics",
"plase":"thrombolytics",
"azepam":"anti-anziety(benzodiazepine)",

```

```
"azine": "antipsychotics (phenothiazine)",
"barbital": "barbiturate",
"dipine": "calcium channel blocker",
"lol": "beta blocker",
"zolam": "cns depressants",
"pril": "ace inhibitor",
"artan": "arb blocker",
"statins": "lipid-lowering drugs",
"parin": "anticoagulants",
"sone": "corticosteroid (prednisone)"}


```

```
def classify_drug(drugname):
    for i in drug_suffix.keys():
        if drugname.endswith(i):
            return drug_suffix[i]


```

```
sensitive['Drug_Class'] =
sensitive['Drug_Term'].apply(classify_drug)


```

```
sensitive


```

```
sensitive['Drug_Class'].unique().tolist()


```

```
sensitive['Drug_Class'].value_counts()


```

```
plt.figure(figsize=(20,10))
sensitive['Drug_Class'].value_counts().plot(kind='bar')
plt.title("Distribution of Drugs By Class")
plt.show()


```

# 🔗 Results

- ◆ Using default spaCy corpus
  - Successful tagging
  - But drugs and illnesses not tagged explicitly

```
Sentence is: Acne is a skin condition caused by dead skin cells sticking together and clogging up pores. Bacteria can play a role, too. A big trigger for the onset of acne is puberty.
/usr/lib/python3.7/rumpy.py:193: UserWarning: [W006] No entities to visualize found in Doc object. If this is surprising to you, make sure the Doc was processed using a model that supports named entity recognition, and che
"__main__", mod_spec)

Acne is a skin condition caused by dead skin cells sticking together and clogging up pores. Bacteria can play a role, too. A big trigger for the onset of acne is puberty.

Sentence is: Doxycycline is a tetracycline antibiotic that fights bacteria in the body. Doxycycline is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, resp
Doxycycline PERSON is a tetracycline antibiotic that fights bacteria in the body. Doxycycline PERSON is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, respiratory infections, eye infections, gonorrhea, chlamydia GPE
, syphilis GPE, periodontitis (gum disease), and others.

Sentence is: Spironolactone, sold under the brand name Aldactone among others, is a medication that is primarily used to treat fluid build-up due to heart failure, liver scarring, or kidney disease.
/usr/lib/python3.7/rumpy.py:193: UserWarning: [W006] No entities to visualize found in Doc object. If this is surprising to you, make sure the Doc was processed using a model that supports named entity recognition, and che
"__main__", mod_spec)

Spironolactone, sold under the brand name Aldactone among others, is a medication that is primarily used to treat fluid build-up due to heart failure, liver scarring, or kidney disease.

Sentence is: Minocycline sold under the brand name Minocin among others, is a tetracycline antibiotic used to treat a number of bacterial infections such as pneumonia.
Minocycline ORG sold under the brand name Minocin PERSON among others, is a tetracycline antibiotic used to treat a number of bacterial infections such as pneumonia.

Sentence is: Aldactone is a prescription medicine used to treat the symptoms of high blood pressure (hypertension), low potassium levels (hypokalemia) and heart failure. Aldactone may be used alone or with other medicatio
/usr/lib/python3.7/rumpy.py:193: UserWarning: [W006] No entities to visualize found in Doc object. If this is surprising to you, make sure the Doc was processed using a model that supports named entity recognition, and che
"__main__", mod_spec)

Aldactone is a prescription medicine used to treat the symptoms of high blood pressure (hypertension), low potassium levels (hypokalemia) and heart failure. Aldactone may be used alone or with other medications.

Sentence is: Clindamycin is an antibiotic used for the treatment of a number of bacterial infections, including bone or joint infections, pelvic inflammatory disease, strep throat, pneumonia, middle ear infections, and en
Clindamycin ORG is an antibiotic used for the treatment of a number of bacterial infections, including bone or joint infections, pelvic inflammatory disease, strep throat PERSON, pneumonia GPE, middle ear GPE infections, and endocarditis.

Sentence is: Tazarotene sold under the brand name Tazorac, among others, is a third-generation prescription topical retinoid.
Tazarotene ORG sold under the brand name Tazorac ORG, among others, is a third ORDINAL -generation prescription topical retinoid.

Sentence is: Tretinoin is a medication used to treat acne and sun-damaged skin. It can't erase deep wrinkles, but it can help improve the appearance of surface wrinkles, fine lines, and darks spots.
/usr/lib/python3.7/rumpy.py:193: UserWarning: [W006] No entities to visualize found in Doc object. If this is surprising to you, make sure the Doc was processed using a model that supports named entity recognition, and che
"__main__", mod_spec)

Tretinoin is a medication used to treat acne and sun-damaged skin. It can't erase deep wrinkles, but it can help improve the appearance of surface wrinkles, fine lines, and darks spots.

Sentence is: Isotretinoin also known as 13-cis-retinoic acid and sold under the brand name Accutane among others, is a medication primarily used to treat severe acne.
Isotretinoin also known as 13-cis CARDINAL -retinoic acid and sold under the brand name Accutane ORG among others, is a medication primarily used to treat severe acne.

Sentence is: Aczone Gel is a prescription medicine used to treat the symptoms of Acne Vulgaris. Aczone Gel may be used alone or with other medications.
Aczone Gel PERSON is a prescription medicine used to treat the symptoms of Acne Vulgaris ORG. Aczone Gel PERSON may be used alone or with other medications.

Sentence is: Adapalene is a third-generation topical retinoid primarily used in the treatment of mild-moderate acne, and is also used off-label to treat keratosis pilaris as well as other skin conditions.
Adapalene is a third ORDINAL -generation topical retinoid primarily used in the treatment of mild-moderate acne, and is also used off-label to treat keratosis pilaris as well as other skin conditions.
```

- ◆ Using the drug and illness corpora
  - modified to show only entities tagged 'DRUG' and 'ILLNESS'
  - Successfully separated drugs and illnesses from data set

```
Sentence is: Acne is a skin condition caused by dead skin cells sticking together and clogging up pores. Bacteria can play a role, too. A big trigger for the onset of acne is puberty.
Acne ILLNESS is a skin condition caused by dead skin cells sticking together and clogging up pores. Bacteria can play a role, too. A big trigger for the onset of acne is puberty.

Sentence is: Doxycycline is a tetracycline antibiotic that fights bacteria in the body. Doxycycline is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, resp
Doxycycline DRUG is a tetracycline DRUG antibiotic that fights bacteria in the body. Doxycycline is used to treat many different bacterial infections, such as acne ILLNESS, urinary tract infections, intestinal infections, respiratory infections, eye infections, gonorrhea, chlamydia,
syphilis, periodontitis (gum disease), and others.

Sentence is: Spironolactone, sold under the brand name Aldactone among others, is a medication that is primarily used to treat fluid build-up due to heart failure, liver scarring, or kidney disease.
Spironolactone DRUG, sold under the brand name Aldactone DRUG among others, is a medication that is primarily used to treat fluid build-up due to heart failure, liver scarring, or kidney disease.

Sentence is: Minocycline sold under the brand name Minocin among others, is a tetracycline antibiotic used to treat a number of bacterial infections such as pneumonia.
Minocycline DRUG sold under the brand name Minocin among others, is a tetracycline DRUG antibiotic used to treat a number of bacterial infections such as pneumonia ILLNESS.

Sentence is: Aldactone is a prescription medicine used to treat the symptoms of high blood pressure (hypertension), low potassium levels (hypokalemia) and heart failure. Aldactone may be used alone or with other medicatio
Aldactone DRUG is a prescription medicine used to treat the symptoms of high blood pressure (hypertension ILLNESS), low potassium levels (hypokalemia) and heart failure. Aldactone may be used alone or with other medications.

Sentence is: Clindamycin is an antibiotic used for the treatment of a number of bacterial infections, including bone or joint infections, pelvic inflammatory disease, strep throat, pneumonia, middle ear infections, and en
Clindamycin DRUG is an antibiotic used for the treatment of a number of bacterial infections, including bone or joint infections, pelvic inflammatory disease, strep throat, pneumonia ILLNESS, middle ear infections, and endocarditis.

Sentence is: Tazarotene sold under the brand name Tazorac, among others, is a third-generation prescription topical retinoid.
Tazarotene DRUG sold under the brand name Tazorac DRUG, among others, is a third-generation prescription topical retinoid.

Sentence is: Tretinoin is a medication used to treat acne and sun-damaged skin. It can't erase deep wrinkles, but it can help improve the appearance of surface wrinkles, fine lines, and darks spots.
Tretinoin DRUG is a medication used to treat acne ILLNESS and sun-damaged skin. It can't erase deep wrinkles, but it can help improve the appearance of surface wrinkles, fine lines, and darks spots.

Sentence is: Isotretinoin also known as 13-cis-retinoic acid and sold under the brand name Accutane among others, is a medication primarily used to treat severe acne.
Isotretinoin DRUG also known as 13-cis-retinoic acid and sold under the brand name Accutane among others, is a medication primarily used to treat severe acne ILLNESS.

Sentence is: Aczone Gel is a prescription medicine used to treat the symptoms of Acne Vulgaris. Aczone Gel may be used alone or with other medications.
Aczone DRUG Gel is a prescription medicine used to treat the symptoms of Acne ILLNESS Vulgaris. Aczone Gel may be used alone or with other medications.

Sentence is: Adapalene is a third-generation topical retinoid primarily used in the treatment of mild-moderate acne, and is also used off-label to treat keratosis pilaris as well as other skin conditions.
Adapalene DRUG is a third-generation topical retinoid primarily used in the treatment of mild-moderate acne ILLNESS, and is also used off-label to treat keratosis pilaris as well as other skin conditions.

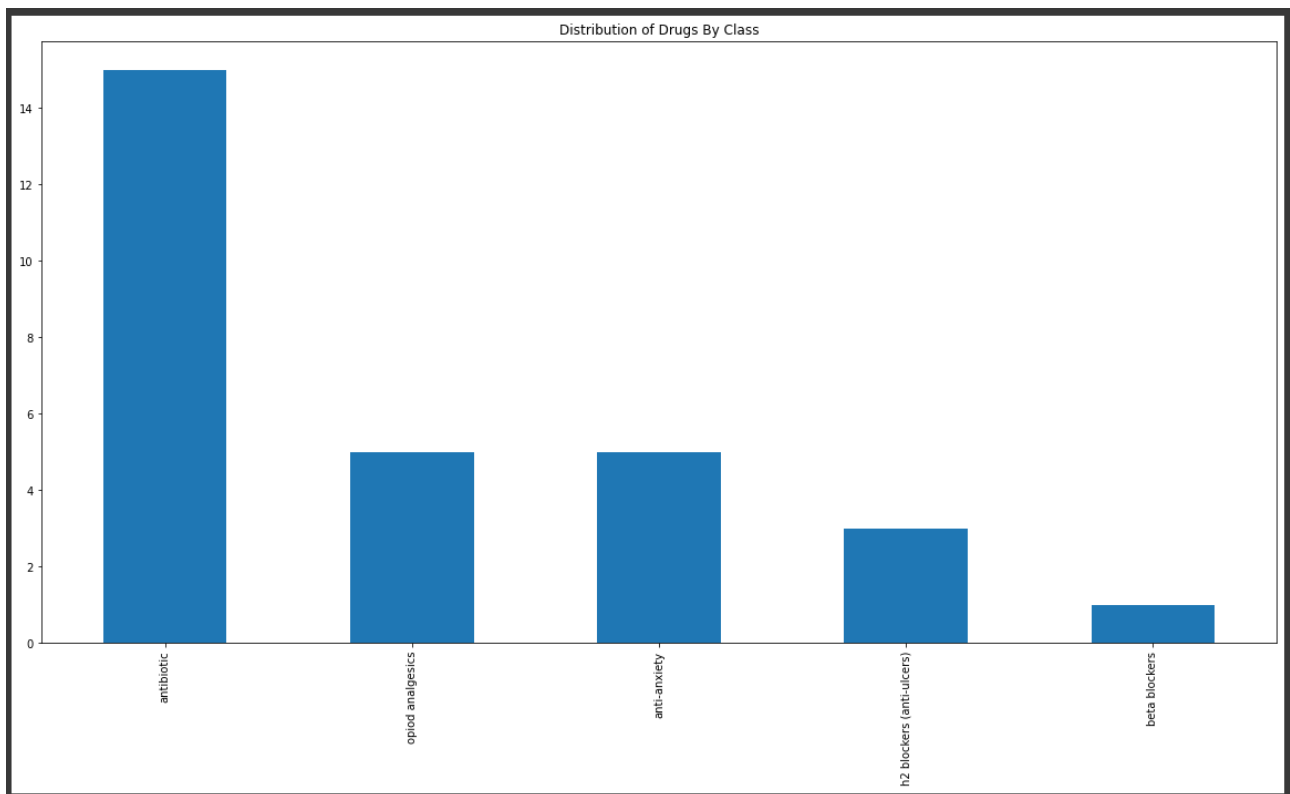
Sentence is: Doryx is a tetracycline antibiotic that fights bacteria in the body. Doryx is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, eye infections,
Doryx DRUG is a tetracycline DRUG antibiotic that fights bacteria in the body. Doryx is used to treat many different bacterial infections, such as acne ILLNESS, urinary tract infections, intestinal infections, eye infections, gonorrhea, chlamydia, periodontitis (gum disease), and
others.

Sentence is: Cefalexin also spelled cephalixin, is an antibiotic that can treat a number of bacterial infections. It kills gram-positive and some gram-negative bacteria by disrupting the growth of the bacterial cell wall.
Cefalexin also spelled cephalixin DRUG, is an antibiotic that can treat a number of bacterial infections. It kills gram-positive and some gram-negative bacteria by disrupting the growth of the bacterial cell wall.

Sentence is: Tetracycline sold under the brand name Sumycin among others, is an oral antibiotic in the tetracyclines family of medications, used to treat a number of infections, including acne, cholera, brucellosis, plagu
Tetracycline DRUG sold under the brand name Sumycin among others, is an oral antibiotic in the tetracyclines family of medications, used to treat a number of infections, including acne ILLNESS, cholera, brucellosis, plague, malaria, and syphilis.

Sentence is: Vibramycin is a tetracycline antibiotic that fights bacteria in the body. Vibramycin is used to treat many different bacterial infections, such as acne, urinary tract infections, intestinal infections, eye in
Vibramycin DRUG is a tetracycline DRUG antibiotic that fights bacteria in the body. Vibramycin is used to treat many different bacterial infections, such as acne ILLNESS, urinary tract infections, intestinal infections, eye infections, gonorrhea, chlamydia, periodontitis (gum
```

◆ Using the drug class to sort using usage statistics





[https://colab.research.google.com/drive/18zX-4MKtloWGu\\_DCL6NEpn5rXTCJoMHQ?usp=sharing](https://colab.research.google.com/drive/18zX-4MKtloWGu_DCL6NEpn5rXTCJoMHQ?usp=sharing)

## Conclusion

- Till now, we have managed to implement acronym and text extraction on our current database.
- We have successfully labelled the DRUG and ILLNESS entities till now, i.e. we can now recognize 22,000+ drug names in the data set and 5000+ illnesses.
- Here NER is used to associate the drug with the illness and is extremely useful for research and analysis
- We have also added multiple statistical methods of data visualization to display the distribution of each drug by class for uses like post market analysis.



## References

- ◆ Text Mining: Applications and Theory - Michael W. Berry Jacob Kogan
- ◆ Text Mining: The state of the art and the challenges - Ah-Hwee Tan
- ◆ What Is Text Mining? - Marti Hearst
- ◆ Biomedical Text Mining and Its Applications - Raul Rodriguez-Esteban  
<https://doi.org/10.1371/journal.pcbi.1000597>
- ◆ A survey of current work in biomedical text mining - Aaron M. Cohen, William R. Hersh  
<https://doi.org/10.1093/bib/6.1.57>
- ◆ <https://www.udemy.com/course/python-for-data-science-and-machine-learning-bootcamp/learn/lecture/5733574#overview>
- ◆  
[https://colab.research.google.com/github/WomenWhoCode/WWCodeDataScience/blob/master/Intro\\_to\\_NLP/2\\_NLP\\_DeepDive1.ipynb](https://colab.research.google.com/github/WomenWhoCode/WWCodeDataScience/blob/master/Intro_to_NLP/2_NLP_DeepDive1.ipynb)
- ◆ <https://archive.ics.uci.edu/ml/datasets/Drug+Review+Dataset+%28Drugs.com%2>
- ◆ <https://www.machinelearningplus.com/spacy-tutorial-nlp/>
- ◆  
[https://github.com/Jcharis/data-science-projects/tree/master/exploratory\\_data\\_analysis\\_in\\_python\\_drug\\_reviews\\_dataset](https://github.com/Jcharis/data-science-projects/tree/master/exploratory_data_analysis_in_python_drug_reviews_dataset)
- ◆, "Named Entity Recognition from Biomedical Text Using SVM," - Z. Ju, J. Wang and F. Zhu  
<https://doi.org/10.1109/icbbe.2011.5779984>
- ◆ A Neural Named Entity Recognition and Multi-Type Normalization Tool for Biomedical Text Mining - Donghyeon Kim, Jinhyuk Lee  
<https://doi.org/10.1109/ACCESS.2019.2920708>