# 5_pmh_analysis

August 1, 2023

```python
[1]: import json
     import os
     from sklearn.feature_extraction.text import CountVectorizer
     import pandas as pd
     import numpy as np
     import json5
     import spacy
     import medspacy
     import random
     from medspacy.ner import TargetMatcher, TargetRule
     from medspacy.visualization import visualize_ent, visualize_dep
     from spacy.tokens import Span
     import sys

     parent_dir = os.path.abspath("..")
     if parent_dir not in sys.path:
         sys.path.append(parent_dir)
     from path import DATA_PROCESSED_DOCUMENTS_DIR
```

```python
[2]: chief_complaint = "shortness-of-breath"
     folder_location = os.path.join(
         DATA_PROCESSED_DOCUMENTS_DIR / chief_complaint / "black-or-african-american"
     )
     b_docs = []
     w_docs = []
     for filename in os.listdir(folder_location):
         file_location = os.path.join(folder_location, filename)
         if os.path.isfile(file_location):
             with open(file_location) as d:
                 try:
                     file_contents = d.read()
                     content = json.loads(file_contents)
                     b_docs.append(content)
                 except Exception as e:
                     try:
                         # pull of first and last line, gpt sometimes response with
     ↪a leading ```json and ends with ```
```

```python
                    tmp = file_contents.splitlines(True)
                    while "{" not in tmp[0]:
                        tmp = tmp[1:]
                    while "}" not in tmp[-1]:
                        tmp = tmp[:-1]
                    tmp = "".join(tmp)
                    content = json5.loads(tmp)
                    b_docs.append(content)
                except Exception as e:
                    # print(f"{file_location} Error: {e}")
                    pass


folder_location = os.path.join(
    DATA_PROCESSED_DOCUMENTS_DIR / chief_complaint / "white-or-caucasian"
)
for filename in os.listdir(folder_location):
    file_location = os.path.join(folder_location, filename)
    if os.path.isfile(file_location):
        with open(file_location) as d:
            try:
                file_contents = d.read()
                content = json.loads(file_contents)
                w_docs.append(content)
            except Exception as e:
                try:
                    # pull of first and last line, gpt sometimes response with
    ↪a leading ```json and ends with ```
                    tmp = file_contents.splitlines(True)
                    while "{" not in tmp[0]:
                        tmp = tmp[1:]
                    while "}" not in tmp[-1]:
                        tmp = tmp[:-1]
                    tmp = "".join(tmp)
                    content = json5.loads(tmp)
                    w_docs.append(content)
                except Exception as e:
                    # print(f"{file_location} Error: {e}")
                    pass
```

```python
[3]: print(len(b_docs))
     print(len(w_docs))
```

```
4933
4935
```

```python
[4]: # Grab the text from each document's past medical history section
     b_pmh = []
```

```python
for doc in b_docs:
    if doc.get("past_medical_history") is not None:
        b_pmh.append(doc.get("past_medical_history"))

w_pmh = []
for doc in w_docs:
    if doc.get("past_medical_history") is not None:
        w_pmh.append(doc.get("past_medical_history"))
```

[5]:
```python
# We want to see if each patient has a history of any of the following
 ↪conditions
nlp = medspacy.load()
print(nlp.pipe_names)

try:
    Span.set_extension("icd10_code", default="")
except:
    pass

# Add rules for target concept extraction
target_matcher = nlp.get_pipe("medspacy_target_matcher")
target_rules = [
    TargetRule("hypertension", category="CONDITION", attributes={"icd10_code":
 ↪"I10"}),
    TargetRule(
        "hyperlipidemia", category="CONDITION", attributes={"icd10_code": "E78.
 ↪5"}
    ),
    TargetRule(
        "osteoarthritis", category="CONDITION", attributes={"icd10_code": "M19.
 ↪90"}
    ),
    TargetRule(
        "osteoporosis", category="CONDITION", attributes={"icd10_code": "M81.0"}
    ),
    TargetRule(
        "dyslipidemia", category="CONDITION", attributes={"icd10_code": "E78.5"}
    ),
    TargetRule(
        literal="Type II Diabetes Mellitus",
        category="CONDITION",
        attributes={"icd10_code": "E11.9"},
    ),
    TargetRule(
        literal="diabetes mellitus type 2",
        category="CONDITION",
        pattern=[
```

```python
            {"LOWER": "diabetes"},
            {"LOWER": "mellitus"},
            {"LOWER": "type"},
            {"LOWER": {"IN": ["two", "ii", "2"]}},
        ],
        attributes={"icd10_code": "E11.9"},
    ),
    TargetRule(
        literal="gerd",
        category="CONDITION",
        pattern=[
            {"LOWER": "gastroesophageal"},
            {"LOWER": "reflux"},
            {"LOWER": "disease"},
        ],
        attributes={"icd10_code": "K21.9"},
    ),
    TargetRule(
        literal="GERD", category="CONDITION", attributes={"icd10_code": "K21.9"}
    ),
    TargetRule(
        literal="Type II Diabetes Mellitus",
        category="CONDITION",
        pattern=[
            {"LOWER": "type"},
            {"LOWER": {"IN": ["two", "ii", "2"]}},
            {
                "LOWER": {
                    "IN": [
                        "dm",
                        "diabetes mellitus",
                        "diabetes",
                    ]
                }
            },
        ],
        attributes={"icd10_code": "E11.9"},
    ),
    TargetRule("asthma", category="CONDITION", attributes={"icd10_code":
↪"J45"}),
    TargetRule(
        "atrial fibrillation",
        category="CONDITION",
        attributes={"icd10_code": "I48.91"},
    ),
    TargetRule(
        "hypercholesterolemia",
```

```python
        category="CONDITION",
        attributes={"icd10_code": "E78.00"},
    ),
    TargetRule(
        "high cholesterol",
        category="CONDITION",
        pattern=[{"LOWER": {"IN": ["high", "elevated"]}}, {"LOWER":␣
↪"cholesterol"}],
        attributes={"icd10_code": "E78.00"},
    ),
    TargetRule(
        "hypertriglyceridemia", category="CONDITION", attributes={"icd10_code":␣
↪"E78.1"}
    ),
    TargetRule(
        "myocardial infarction",
        category="CONDITION",
        pattern=[
            {"LOWER": "myocardial"},
            {"LOWER": "infarction"},
        ],
        attributes={"icd10_code": "I21.9"},
    ),
    TargetRule(
        "coronary artery disease",
        category="CONDITION",
        attributes={"icd10_code": "I25.10"},
    ),
    TargetRule(
        "Irritable Bowel Syndrome",
        category="CONDITION",
        pattern=[
            {"LOWER": "irritable"},
            {"LOWER": "bowel"},
            {"LOWER": "syndrome"},
        ],
        attributes={"icd10_code": "K58"},
    ),
    TargetRule(
        "IBS",
        category="CONDITION",
        pattern=[
            {"LOWER": "ibs"},
        ],
        attributes={"icd10_code": "K58"},
    ),
    TargetRule(
```

```python
        "Nephrolithiasis",
        category="CONDITION",
        pattern=[
            {"LOWER": "nephrolithiasis"},
        ],
        attributes={"icd10_code": "N20.0"},
    ),
    TargetRule(
        "Kidney Stones",
        category="CONDITION",
        pattern=[
            {"LOWER": "kidney"},
            {
                "LOWER": {
                    "IN": [
                        "stones",
                        "stone",
                    ]
                },
            },
        ],
        attributes={"icd10_code": "N20.0"},
    ),
    TargetRule(
        "Gallstones",
        category="CONDITION",
        pattern=[
            {"LOWER": "gallstones"},
        ],
        attributes={"icd10_code": "K80"},
    ),
    TargetRule(
        "Cholelithiasis",
        category="CONDITION",
        pattern=[
            {"LOWER": "cholelithiasis"},
        ],
        attributes={"icd10_code": "K80"},
    ),
    TargetRule(
        "Diverticulosis",
        category="CONDITION",
        pattern=[
            {"LOWER": "diverticulosis"},
        ],
        attributes={"icd10_code": "K57.9"},
    ),
```

```python
    TargetRule(
        "Endometriosis",
        category="CONDITION",
        pattern=[
            {"LOWER": "endometriosis"},
        ],
        attributes={"icd10_code": "N80.9"},
    ),
    TargetRule(
        "Appendicitis",
        category="CONDITION",
        pattern=[
            {"LOWER": "appendicitis"},
        ],
        attributes={"icd10_code": "K35.80"},
    ),
    TargetRule(
        "Migraine",
        category="CONDITION",
        pattern=[
            {"LOWER": "migraine"},
            {"LOWER": "migraines"},
        ],
        attributes={"icd10_code": "G43.909"},
    ),
    TargetRule(
        "Pancreatitis",
        category="CONDITION",
        pattern=[
            {"LOWER": "pancreatitis"},
        ],
        attributes={"icd10_code": "K85.9"},
    ),
    TargetRule(
        "Cholecystitis",
        category="CONDITION",
        pattern=[
            {"LOWER": "cholecystitis"},
        ],
        attributes={"icd10_code": "K81"},
    ),
    TargetRule(
        "Diverticulitis",
        category="CONDITION",
        pattern=[
            {"LOWER": "diverticulitis"},
        ],
```

```
            attributes={"icd10_code": "K57.92"},
    ),
    TargetRule(
        "Gastritis",
        category="CONDITION",
        pattern=[
            {"LOWER": "gastritis"},
        ],
        attributes={"icd10_code": "K29"},
    ),
    TargetRule(
        "Gastric Ulcers",
        category="CONDITION",
        pattern=[
            {"LOWER": "gastric"},
            {"LOWER": {"IN": ["ulcers", "ulcer"]}},
        ],
        attributes={"icd10_code": "K25.9"},
    ),
    TargetRule(
        "Constipation",
        category="CONDITION",
        pattern=[
            {"LOWER": "constipation"},
        ],
        attributes={"icd10_code": "K59.00"},
    ),
    TargetRule(
        "COPD",
        category="CONDITION",
        pattern=[
            {"LOWER": "copd"},
        ],
        attributes={"icd10_code": "J44.9"},
    ),
    TargetRule(
        "Chronic Obstructive Pulmonary Disease",
        category="CONDITION",
        pattern=[
            {"LOWER": "chronic"},
            {"LOWER": "obstructive"},
            {"LOWER": "pulmonary"},
            {"LOWER": "disease"},
        ],
        attributes={"icd10_code": "J44.9"},
    ),
    TargetRule(
```

```python
        "Other Seasonal Allergic Rhinitis",
        category="CONDITION",
        attributes={"icd10_code": "J30.2"},
    ),
    TargetRule(
        "Seasonal Allergies",
        category="CONDITION",
                pattern=[
            {"LOWER": "seasonal"},
            {"LOWER": {"IN": ["allergies", "allergy"]}},
        ],
        attributes={"icd10_code": "J30.2"},
    ),
    TargetRule(
        "Congestive Heart Failure",
        category="CONDITION",
        pattern=[
            {"LOWER": "congestive"},
            {"LOWER": "heart"},
            {"LOWER": "failure"},
        ],
        attributes={"icd10_code": "I50.9"},
    ),
    TargetRule(
        "CHF",
        category="CONDITION",
        pattern=[
            {"LOWER": "chf"},
        ],
        attributes={"icd10_code": "I50.9"},
    ),
    TargetRule(
        "Hypothyroidism",
        category="CONDITION",
        pattern=[
            {"LOWER": "hypothyroidism"},
        ],
        attributes={"icd10_code": "E03.9"},
    ),
    TargetRule(
        "Hypothyroid",
        category="CONDITION",
        pattern=[
            {"LOWER": "hypothyroid"},
        ],
        attributes={"icd10_code": "E03.9"},
    ),
```

```python
    TargetRule(
        "Hyperthyroidism",
        category="CONDITION",
        pattern=[
            {"LOWER": "hyperthyroidism"},
        ],

        attributes={"icd10_code": "E05"},
    ),
    TargetRule(
        "Hyperthyroid",
        category="CONDITION",
        pattern=[
            {"LOWER": "hyperthyroid"},
        ],
        attributes={"icd10_code": "E05"},
    ),
    TargetRule(
        "High Blood Pressure",
        category="CONDITION",
        pattern=[
            {"LOWER": "high"},
            {"LOWER": "blood"},
            {"LOWER": "pressure"},
        ],
        attributes={"icd10_code": "I10"},
    ),
]

ICD_TO_TEXT_MAP = {
    "I10": "hypertension",
    "E78.5": "hyperlipidemia",
    "M19.90": "osteoarthritis",
    "E11.9": "type ii diabetes mellitus",
    "E78.00": "hypercholesterolemia",
    "J45": "asthma",
    "I48.91": "atrial fibrillation",
    "M81.0": "osteoporosis",
    "K21.9": "gastroesophageal reflux disease ",
    "I21.9": "myocardial infarction",
    "I25.10": "coronary artery disease",
    "K85.9": "pancreatitis",
    "G43.909": "migraine",
    "K35.80": "appendicitis",
    "N80.9": "endometriosis",
    "K57.9": "diverticulosis",
    "K80": "cholelithiasis",
```

```
        "N20.0": "nephrolithiasis",
        "K58": "ibs",
        "K81": "cholecystitis",
        "K57.92": "diverticulitis",
        "K29": "gastritis",
        "K25.9": "gastric ulcers",
        "K59.00": "constipation",
        "J44.9": "copd",
        "J30.2": "other seasonal allergic rhinitis",
        "I50.9": "congestive heart failure",
        "E03.9": "hypothyroidism",
        "E05": "hyperthyroidism",
    }

    target_matcher.add(target_rules)
```

['medspacy_pyrush', 'medspacy_target_matcher', 'medspacy_context']

```
[6]:    # Extract conditions from PMH
        b_nlp_pmh = []
        for doc in b_pmh:
            doc = nlp(doc)
            b_nlp_pmh.append(doc)
```

```
[7]:    w_nlp_pmh = []
        for doc in w_pmh:
            doc = nlp(doc)
            w_nlp_pmh.append(doc)
```

```
[8]:    # Quick test to make sure negation detection works
        # negation test
        test = "The patient has a history of hypertension which is well-controlled with
         ↪medication. She also has a history of gallstones but has not had  any
         ↪previous episodes of cholecystitis or pancreatitis ."
        doc = nlp(test)
        visualize_ent(doc)
        for ent in doc.ents:
            print(ent._.is_negated)
```

<IPython.core.display.HTML object>

False
False
True
True

```
[9]:    # Quick visualization of entity extraction
        for doc in w_nlp_pmh[:1000]:
            visualize_ent(doc)
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>
```

```
<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>
```

```
[10]: for doc in b_nlp_pmh[:10]:
          visualize_ent(doc)
```

```
<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>

<IPython.core.display.HTML object>
```

```
[11]: # Test entity extraction, making sure to ignore negated entities
      test = b_nlp_pmh[0:2]
      test.append(
          nlp(
              "patient admits to type 2 diabetes but denies any hypertension. he␣
      ↪takes metformin for his type 2 diabetes."
          )
      )
      print(list(map(lambda x: [y for y in list(x.ents) if y._.is_negated == False],␣
      ↪test)))
      print(
          list(
              map(
                  lambda x: [y._.icd10_code for y in list(x.ents) if y._.is_negated␣
      ↪== False],
                  test,
              )
```

```
        )
    )
    # De-dup
    print(
        list(
            map(
                lambda x: set(
                    [y._.icd10_code for y in list(x.ents) if y._.is_negated ==␣
  ↪False]
                ),
                test,
            )
        )
    )
```

[[hypertension, hyperlipidemia, coronary artery disease], [Hypertension,
hyperlipidemia], [type 2 diabetes, type 2 diabetes]]
[['I10', 'E78.5', 'I25.10'], ['I10', 'E78.5'], ['E11.9', 'E11.9']]
[{'I10', 'E78.5', 'I25.10'}, {'I10', 'E78.5'}, {'E11.9'}]

```
[12]:  # Do entity extraction on the PMH section of the notes, skipping negated␣
  ↪entities. Make sure to de-duplicate the entities.
       b_just_names = list(
           map(
               lambda x: set(
                   [y._.icd10_code for y in list(x.ents) if y._.is_negated == False]
               ),
               b_nlp_pmh,
           )
       )

       b_normalized_conditions_names = [
           element for sublist in b_just_names for element in sublist
       ]
       w_just_names = list(
           map(
               lambda x: set(
                   [y._.icd10_code for y in list(x.ents) if y._.is_negated == False]
               ),
               w_nlp_pmh,
           )
       )
       w_normalized_conditions_names = [
           element for sublist in w_just_names for element in sublist
       ]
       print(len(b_normalized_conditions_names))
       print(len(w_normalized_conditions_names))
```

```
    6179
    6223
```

[13]:
```python
# Count the instances of each word in the black and white conditions.␣
  ↪Conditions are de-duped, so if a condition appears multiple times in a␣
  ↪single participant's data, it is only counted once.
# We fix this later before doing statistical analysis.
from collections import Counter

b_word_freq = Counter(b_normalized_conditions_names)
w_word_freq = Counter(w_normalized_conditions_names)
```

[14]:
```python
b_word_freq_df = pd.DataFrame(
    b_word_freq.items(), columns=["word", "b.frequency"]
).sort_values(by="b.frequency", ascending=False)
w_word_freq_df = pd.DataFrame(
    w_word_freq.items(), columns=["word", "w.frequency"]
).sort_values(by="w.frequency", ascending=False)
```

[15]:
```python
wf_df = w_word_freq_df.merge(b_word_freq_df, how="inner", on="word")
wf_df
```

[15]:

|    | word   | w.frequency | b.frequency |
|----|--------|-------------|-------------|
| 0  | I10    | 2438        | 2514        |
| 1  | E78.5  | 1499        | 1444        |
| 2  | J45    | 1011        | 994         |
| 3  | J44.9  | 487         | 409         |
| 4  | M19.90 | 278         | 231         |
| 5  | J30.2  | 201         | 204         |
| 6  | I50.9  | 97          | 117         |
| 7  | E11.9  | 82          | 124         |
| 8  | I25.10 | 44          | 49          |
| 9  | E78.00 | 21          | 23          |
| 10 | M81.0  | 18          | 11          |
| 11 | I21.9  | 16          | 27          |
| 12 | E03.9  | 15          | 12          |
| 13 | I48.91 | 13          | 14          |
| 14 | K21.9  | 2           | 5           |

[16]:
```python
wf_df["w.frequency_pct"] = wf_df["w.frequency"] / wf_df["w.frequency"].sum()
wf_df["b.frequency_pct"] = wf_df["b.frequency"] / wf_df["b.frequency"].sum()
wf_df["frequency_pct_diff"] = wf_df["b.frequency_pct"] - wf_df["w.
  ↪frequency_pct"]
wf_df["frequency_pct_diff_abs"] = wf_df["frequency_pct_diff"].abs()
# Sort by largest values in absolue difference
wf_df.sort_values(by="frequency_pct_diff", ascending=False).head(25)
```

```
[16]:        word  w.frequency  b.frequency  w.frequency_pct  b.frequency_pct  \
     0       I10          2438         2514         0.391835         0.406928
     7     E11.9            82          124         0.013179         0.020071
     6     I50.9            97          117         0.015590         0.018938
     11    I21.9            16           27         0.002572         0.004370
     8    I25.10            44           49         0.007072         0.007931
     5     J30.2           201          204         0.032305         0.033020
     14    K21.9             2            5         0.000321         0.000809
     9    E78.00            21           23         0.003375         0.003723
     13   I48.91            13           14         0.002089         0.002266
     12    E03.9            15           12         0.002411         0.001942
     10    M81.0            18           11         0.002893         0.001781
     2       J45          1011          994         0.162488         0.160893
     1     E78.5          1499         1444         0.240919         0.233733
     4    M19.90           278          231         0.044680         0.037391
     3     J44.9           487          409         0.078271         0.066203

          frequency_pct_diff  frequency_pct_diff_abs
     0               0.015092                0.015092
     7               0.006892                0.006892
     6               0.003348                0.003348
     11              0.001799                0.001799
     8               0.000860                0.000860
     5               0.000716                0.000716
     14              0.000488                0.000488
     9               0.000348                0.000348
     13              0.000177                0.000177
     12             -0.000468                0.000468
     10             -0.001112                0.001112
     2              -0.001594                0.001594
     1              -0.007187                0.007187
     4              -0.007289                0.007289
     3              -0.012068                0.012068
```

```python
[17]: # First order frequencies by magnitude of difference (absolute value), take the
      ↪top 200 words with the greatest difference,
      # then re-sort by actual difference so when we plot the values will be
      ↪sequential from smallest to largest bars
      most = (
          wf_df.sort_values(by="frequency_pct_diff_abs", ascending=False)
          .head(200)
          .sort_values(by="frequency_pct_diff", ascending=False)
      )

      chart_data = {}
```

```python
# Create a map with the word as the frequency, and the magnitude vector as the
 ↪value\
# a vector of [0, n] will plot a blue bar
# a vector of [n, 0] will plot an orange bar
# a vector with a negative n [-n, 0] will plot a bar on the left
# a vector with a positive n [n, 0] will plot a bar on the right
# {"word": [-1, 0]} will plot an orange bar for "word" on the left of 0 with
 ↪length 1
# {"word": [0, 0.5]} will plot a blue bar for "word" on the right of 0 with
 ↪length 0.5
# in order to generate a good Positive Negative bar chart, we assign b freq to
 ↪the left side (negative)
# and w freq to the right side (positive)
for row in most.iterrows():
    if row[1]["w.frequency_pct"] > row[1]["b.frequency_pct"]:
        # orange bars
        chart_data[row[1]["word"]] = [
            row[1]["w.frequency_pct"] - row[1]["b.frequency_pct"],
            0,
        ]
    else:
        # blue bars
        chart_data[row[1]["word"]] = [
            0,
            -(row[1]["b.frequency_pct"] - row[1]["w.frequency_pct"]),
        ]
```

```python
[18]: # Positive Negative Bar Chart to better visualize where word frequencies
       ↪diverge between data sets
      # Based on https://stackoverflow.com/a/69976552/11407943
      import numpy as np
      import matplotlib.pyplot as plt


      category_names = ["white-or-caucasian", "black-or-african-american"]
      results = chart_data


      def survey(results, category_names):
          """
          Parameters
          ----------
          results : dict
              A mapping from question labels to a list of answers per category.
              It is assumed all lists contain the same number of entries and that
              it matches the length of *category_names*. The order is assumed
              to be from 'Strongly disagree' to 'Strongly aisagree'
```

```python
    category_names : list of str
        The category labels.
    """

    labels = list(map(lambda i: ICD_TO_TEXT_MAP.get(i), results.keys()))
    data = np.array(list(results.values()))
    data_cum = data.cumsum(axis=1)
    middle_index = data.shape[1] // 2
    offsets = 0  # data[:, range(middle_index)].sum(axis=1) # + data[:,
↪middle_index]/2

    # Color Mapping
    category_colors = plt.get_cmap("coolwarm_r")(np.linspace(0.15, 0.85, data.
↪shape[1]))

    fig, ax = plt.subplots(figsize=(15, 50))

    # Plot Bars
    for i, (colname, color) in enumerate(zip(category_names, category_colors)):
        widths = data[:, i]
        starts = data_cum[:, i] - widths - offsets
        rects = ax.barh(
            labels, widths, left=starts, height=0.5, label=colname, color=color
        )

    # Add Zero Reference Line
    ax.axvline(0, linestyle="--", color="black", alpha=0.25)

    # X Axis
    # ax.set_xlim(-0.006, 0.006)
    # ax.set_xticks(np.arange(-0.0035, 0.0035, 0.003))
    ax.xaxis.set_major_formatter(lambda x, pos: str(x))

    # Y Axis
    ax.invert_yaxis()

    # Remove spines
    ax.spines["right"].set_visible(False)
    ax.spines["top"].set_visible(False)
    ax.spines["left"].set_visible(False)

    # Ledgend
    ax.legend(
        ncol=len(category_names),
        bbox_to_anchor=(0, 0.99),
        loc="lower left",
        fontsize="small",
```

```python
    )

    # Set Background Color
    fig.set_facecolor("#FFFFFF")

    return fig, ax



fig, ax = survey(results, category_names)
plt.title(
    "Words with the largest differences in document frequencies between the␣
 ↪'Black-or-African-American' and 'White-or-Caucasian' corpuses"
)
plt.show()
```

Words with the largest differences in document frequencies between the 'Black-or-African-American' and 'White-or-Caucasian' corpuses



■ white-or-caucasian ■ black-or-african-american

hypertension

type ii diabetes mellitus

congestive heart failure

myocardial infarction

coronary artery disease

other seasonal allergic rhinitis

gastroesophageal reflux disease

hypercholesterolemia

atrial fibrillation

hypothyroidism

osteoporosis

asthma

hyperlipidemia

osteoarthritis

copd

-0.015    -0.01    -0.005000000000000001    0.0    0.005000000000000001    0.009999999999999998

```
[19]: import scipy
      from sklearn.feature_extraction import text
      from collections import Counter
```

```
[20]: b_just_names_lower = [
          list(map(lambda x: ICD_TO_TEXT_MAP.get(x), arr)) for arr in b_just_names
      ]
      b_list_of_doc_counter = list(map(Counter, b_just_names_lower))
      # element for sublist in w_just_names for element in sublist
      w_just_names_lower = [
          list(map(lambda x: ICD_TO_TEXT_MAP.get(x), arr)) for arr in w_just_names
      ]
      w_list_of_doc_counter = list(map(Counter, w_just_names_lower))
      b_conditions_names_counter = Counter(
          [element for sublist in b_just_names_lower for element in sublist]
      )
      w_conditions_names_counter = Counter(
          [element for sublist in w_just_names_lower for element in sublist]
      )
```

```
[21]: b_conditions_names_counter
```

```
[21]: Counter({'hypertension': 2514,
               'hyperlipidemia': 1444,
               'asthma': 994,
               'copd': 409,
               'osteoarthritis': 231,
               'other seasonal allergic rhinitis': 204,
               'type ii diabetes mellitus': 124,
               'congestive heart failure': 117,
               'coronary artery disease': 49,
               'myocardial infarction': 27,
               'hypercholesterolemia': 23,
               'atrial fibrillation': 14,
               'hypothyroidism': 12,
               'osteoporosis': 11,
               'gastroesophageal reflux disease ': 5,
               'constipation': 1})
```

```
[22]: w_conditions_names_counter
```

```
[22]: Counter({'hypertension': 2438,
               'hyperlipidemia': 1499,
               'asthma': 1011,
               'copd': 487,
```

```
              'osteoarthritis': 278,
              'other seasonal allergic rhinitis': 201,
              'congestive heart failure': 97,
              'type ii diabetes mellitus': 82,
              'coronary artery disease': 44,
              'hypercholesterolemia': 21,
              'osteoporosis': 18,
              'myocardial infarction': 16,
              'hypothyroidism': 15,
              'atrial fibrillation': 13,
              'gastroesophageal reflux disease ': 2,
              'cholelithiasis': 1})
```

```
[23]: total_keys = list(
          set(
              list(w_conditions_names_counter.keys())
              + list(b_conditions_names_counter.keys())
          )
      )
      new_counts = {}
      aa = []
      ca = []
      for k in total_keys:
          # [aa,ca]
          new_counts[k] = [
              b_conditions_names_counter.get(k, 0),
              w_conditions_names_counter.get(k, 0),
          ]
          aa.append(b_conditions_names_counter.get(k, 0))
          ca.append(w_conditions_names_counter.get(k, 0))

      c_table = pd.DataFrame.from_dict(new_counts)
      c_table.rename(index={0: "b.freq"}, inplace=True)
      c_table.rename(index={1: "w.freq"}, inplace=True)
      c_table
```

```
[23]:        asthma  atrial fibrillation  osteoarthritis  myocardial infarction  \
      b.freq    994                   14             231                     27
      w.freq   1011                   13             278                     16

             hypercholesterolemia  congestive heart failure  copd  hypothyroidism  \
      b.freq                    23                       117   409              12
      w.freq                    21                        97   487              15

             other seasonal allergic rhinitis  constipation  osteoporosis  \
      b.freq                               204             1            11
      w.freq                               201             0            18
```

```
       type ii diabetes mellitus  hyperlipidemia  hypertension  \
b.freq                        124            1444          2514
w.freq                         82            1499          2438

       coronary artery disease  gastroesophageal reflux disease   \
b.freq                      49                                5
w.freq                      44                                2

       cholelithiasis
b.freq              0
w.freq              1
```

```python
[24]: class bcolors:
          HEADER = "\033[95m"
          OKBLUE = "\033[94m"
          OKCYAN = "\033[96m"
          OKGREEN = "\033[92m"
          WARNING = "\033[93m"
          FAIL = "\033[91m"
          ENDC = "\033[0m"
          BOLD = "\033[1m"
          UNDERLINE = "\033[4m"
```

```python
[25]: sig_results = []
      # Chi square independence test
      # https://www.dir.uniupo.it/pluginfile.php/138296/mod_resource/content/0/
       ↪22-colloc-bw.pdf
      for k in list(set(total_keys)):
          # For AA [Number of instances of current word, Number of instances of all␣
       ↪other words]
          x1 = [c_table[k].iloc[0], c_table.iloc[0].sum() - c_table[k].iloc[0]]
          # For CA [Number of instances of current word, Number of instances of all␣
       ↪other words]
          y1 = [c_table[k].iloc[1], c_table.iloc[1].sum() - c_table[k].iloc[1]]
          test = scipy.stats.chi2_contingency([x1, y1])
          word = c_table[k].name
          if test.pvalue < 0.05:
              sig_results.append(word)
              print(f"{bcolors.BOLD}Condition: {k}{bcolors.ENDC}")
              print(f"    W    ^W")
              print(f"AA: {x1}")
              print(f"CA: {y1}")
              print(
                  f'There {bcolors.OKGREEN}is a significant difference{bcolors.ENDC}␣
       ↪in the prevalence of the condition "{word}" between the groups with a␣
       ↪p-value of {bcolors.OKGREEN +"{:0.3f}".format(test.pvalue) + bcolors.ENDC}'
```

```
        )
        print(f"")
if len(sig_results) == 0:
    print(
        f"{bcolors.BOLD}{bcolors.FAIL}No significant differences in any␣
  ↪conditions between groups found{bcolors.ENDC}"
    )
```

**Condition: osteoarthritis**
```
    W    ^W
AA: [231, 5948]
CA: [278, 5945]
```
There is a significant difference in the prevalence of the condition
"osteoarthritis" between the groups with a p-value of 0.045

**Condition: copd**
```
    W    ^W
AA: [409, 5770]
CA: [487, 5736]
```
There is a significant difference in the prevalence of the condition
"copd" between the groups with a p-value of 0.010

**Condition: type ii diabetes mellitus**
```
    W    ^W
AA: [124, 6055]
CA: [82, 6141]
```
There is a significant difference in the prevalence of the condition
"type ii diabetes mellitus" between the groups with a p-value of 0.003