

In []:

```
'''
import xml.etree.ElementTree
e = xml.etree.ElementTree.parse(url)

from bs4 import BeautifulSoup

y=BeautifulSoup(e)

import requests
import xml.etree.ElementTree as ET

r = requests.get(url)
root = ET.fromstring(r.text)

#from bs4 import BeautifulSoup
#y=BeautifulSoup(r)

print (r)


import json

from lxml import etree

import requests
import xml.etree.ElementTree as ET

r = requests.get(url)
root = ET.fromstring(r.text)

dom = etree.parse(r)
# Load XSLT
transform = etree.XSLT(etree.fromstring(XSL))

# apply XSLT on Loaded dom
json_text = str(transform(dom))

# json_text contains the data converted to JSON format.
# you can use it with the JSON API. Example:
data = json.loads(json_text)
print(data)

'''

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings"
```

```

rings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional
_hearings_from_server/gpo_tools/metadata_results_new.csv"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-105hhrg40050']
count = 0
for jacket in df1['filename']:

    # try:
        print (count)
        url = 'https://api.govinfo.gov/packages/'+jacket+'/mods?&api_key=XNEgGxjbEszIMy
Ieni9xpgdkqy60QD5p9S4Vvdlc'

        r = requests.get(url)

        with open('data.xml', 'w') as f:
            f.write(r.text)

        with open("data.xml", 'r') as f:
            xmlString = f.read()

        #print ("XML input (data.xml):")
        #print(xmlString)

        jsonString = json.dumps(xmltodict.parse(xmlString), indent=4)

        jsonObj = json.loads(jsonString)

        #print("\nJSON output(output.json):")
        #print(jsonString)

        #with open("output.json", 'w') as f:
        #    f.write(jsonString)

        witnesses = []
        witness_count = 0
        try:
            if "witness" in jsonObj["mods"]["extension"][2]:
                for witness in (jsonObj["mods"]["extension"][2]["witness"]):
                    witnesses.append(witness+'\n')
                    witness_count += 1
        except:
            witnesses.append ("Not found\n")

        count = count + 1

    print ("".join(witnesses))

    with open(metadata_results, 'r') as csvinput:
        with open(metadata_results_new, 'a') as csvoutput:
            writer = csv.writer(csvoutput, lineterminator='\n')
            reader = csv.reader(csvinput)

```

```
all = []
row = next(reader)
row.append('Witnesses & Affiliattions')
all.append(row)

for row in reader:
    row.append("".join(witnesses))
    all.append(row)

writer.writerows(all)

if (count > 2):
    break

#except:
#    count = count + 1
#    continue
```

Congressional committee name:

In []:

```
print (jsonObj["mods"]["name"][0]["namePart"])
```

Witnesses:

In []:

```
witness_count = 0
if "witness" in jsonObj["mods"]["extension"][2]:
    for witness in (jsonObj["mods"]["extension"][2]["witness"]):
        print (witness)
        witness_count += 1
else:
    print ("No witness information found")
```

Affiliations:

In []:

```
nameAff = {}
for name in (jsonObj["mods"]["name"]):
    if name["@type"] == "personal" and "affiliation" in name:
        nameAff[name['namePart']] = name['affiliation']

for i in nameAff.items():
    print (i[0] + '\t' + i[1])
```

In []:

```
# Metadata_results
# Committee number column - from individual csv

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csv = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"

committees = {}

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = [ 'CHRG-115hhrg27211' ]
count = 0
for jacket in df1['filename']:

    try:
        #print (count)

        #if (count > 50):
        #    break

        count = count + 1

        df2 = pd.read_csv(results_csv+jacket+'.csv')

        committees[jacket] = df2['committees'].iloc[0]

    except:
        count = count + 1
        continue

print (committees)
```

In []:

```

# Metadata_results
# Committee number column - from individual csv

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csv = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = [ 'CHRG-115hrg27211' ]
count = 0

with open(metadata_results, 'r') as csvinput:
    with open(metadata_results_new, 'w') as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)
        row.append('Committees')
        all.append(row)

        for row in reader:

            try:

                if ( not math.isnan(committees[row[5]]) ):
                    row.append(committees[row[5]])

            else:
                row.append("-")
            except:
                row.append("-")

            all.append(row)

        writer.writerows(all)

```

In []:

```
# Individual CSVs
# Affiliations

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hrg27211']
```

```

count = 0

files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

for file in files:

    try:

        url = 'https://api.govinfo.gov/packages/'+file.strip()[:-4]+' /mods?&api_key=qv5
08dpECfRcX6wttIoMw63RT81NPRgkNpsU58c2'

        #print (url)
        r = requests.get(url)

        with open('data.xml', 'w', encoding="utf8") as f:
            f.write(r.text)

        with open("data.xml", 'r', encoding="utf8") as f:
            xmlString = f.read()

        #print ("XML input (data.xml):")
        #print(xmlString)

        jsonString = json.dumps(xmltodict.parse(xmlString), indent=4)
        jsonObj = json.loads(jsonString)

        with open(results_csvs+file, 'r', encoding="utf8") as csvinput:
            with open(results_csvs_new+file, 'w+', encoding="utf8") as csvoutput:
                writer = csv.writer(csvoutput, lineterminator='\n')
                reader = csv.reader(csvinput)

                all = []
                row = next(reader)

                row.append('Full name')
                row.append('Affiliation')
                all.append(row)
                #print (row)
                #try:

                for row in reader:

                    try:

                        if ( row[-1] == "Yes"):
                            row.append("".join(row[5].split(",")[:2]).strip())
                            row.append("".join(row[5].split(",")[:2]).strip())
                        else:

                            try:
                                nameAff = {}
                                for name in (jsonObj["mods"]["name"]):
                                    if name["@type"] == "personal" and "affiliatio
n" in name:
                                        nameAff[name['namePart']] = name['affiliati
on']

                                added = False
                                for i in nameAff.items():
                                    if (fuzz.token_sort_ratio(i[0], row[5].strip())

```



```
> 85):  
  
        row.append(i[0])  
        row.append(i[1])  
        added = True  
        break  
  
    if(not added):  
        row.append(row[5].strip())  
        row.append("-")  
  
    except:  
        row.append(row[5].strip())  
        row.append("-")  
  
    except:  
        row.append(row[5].strip())  
        row.append("-")  
  
    all.append(row)  
  
    #except:  
    #    writer.writerow(all)  
    #    continue  
    writer.writerow(all)  
  
except:  
    continue
```

In []:

```

# Metadata_results
# Witness names & Affiliations, Members of the congress

import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

```

```

all = []
row = next(reader)

row.append('Witnesses')
row.append('Members of the congress')
row.append('File exists')
all.append(row)
#print (row)
#try:

for row in reader:

    #try:
    if (row[6].strip()+'.csv' in os.listdir(results_csvs)):
        print (row[6].strip()+'.csv')
        file = pd.read_csv(results_csvs + row[6].strip() + '.cs
v')

        #print (file.head())
        witnesses = []
        members = []

        for index, row1 in file.iterrows():
            #print (row1['Witness'])
            temp = ''
            if (row1['Witness'].strip() == "Yes"):
                if (str(row1['Full name']).strip() != 'NA' and
str(row1['Full name']).strip() != '-' and str(row1['Full name']).strip() != ''):
                    temp = str(row1['Full name'])
                    if (str(row1['Affiliation']).strip() != 'N
A' and str(row1['Affiliation']).strip() != '-' and str(row1['Affiliation']).strip() !=
''):
                        temp += ' : ' + str(row1['Affiliation'
]).strip() + ';\n'

                    witnesses.append(temp)
                else:
                    witnesses.append(temp + ';\n')
            else:
                if (str(row1['Full name']).strip() != 'NA' and
str(row1['Full name']).strip() != '-' and str(row1['Full name']).strip() != ''):
                    temp = str(row1['Full name'])
                    if (str(row1['Affiliation']).strip() != 'N
A' and str(row1['Affiliation']).strip() != '-' and str(row1['Affiliation']).strip() !=
''):
                        temp += ' : ' + str(row1['Affiliation'
]).strip() + ';\n'

                    members.append(temp)
                else:
                    members.append(temp + ';\n')

            #print (witnesses)

witnesses = [x for x in witnesses if str(x) != 'nan;']
members = [x for x in members if str(x) != 'nan;']

witnesses = set(witnesses)
members = set(members)

if (len(witnesses) == 0):

```

```
        row.append(' - ')
    else:
        row.append("".join(witnesses).strip())

    if (len(members) == 0):
        row.append(' - ')
    else:
        row.append("".join(members).strip())

    row.append("Yes")

    all.append(row)

else:
    row.append(' - ')
    row.append(' - ')
    row.append("No")
    all.append(row)

except:
    # row.append("- ")
    # row.append("- ")

    # all.append(row)
    # continue
except:
    # writer.writerow(all)
    # continue
writer.writerow(all)
```

In []:

```

# GPO agencies
# Individual CSVs

import csv
import pandas as pd
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

df = pd.read_csv(gpo)
agencies = []
for i in df['Agency']:
    temp = i.replace('U.S.', 'United States')
    temp = temp.replace('U.S', 'United States')
    temp = temp.replace('Dep.', 'Department')

    agencies.append(temp)

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

file = pd.read_csv(sample_csvs + 'CHRG-104hhrg37344' + '.csv')
for index, row1 in file.iterrows():
    if (row1['Witness'] == "Yes"):
        max_score = 0
        for i in set(agencies):
            score = fuzz.token_set_ratio( i.lower(), row1['Affiliation'].lower())
            if (score > max_score):
                max_score = score
                agency = i
        print ( row1['Affiliation'] + ' : ' + agency + '\t' + str(max_score))

```

In []:

```

# metadata_results_new
# Remove "nan"

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csv = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = [ 'CHRG-115hhrg27211' ]
count = 0

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        all.append(row)

        for row in reader:

            row[-2] = "\n".join( list(filter(None, row[-2].replace('nan;', '').split("\n"))))

            if(row[-2].strip() == ''):
                row[-2] = '-'

            all.append(row)

        writer.writerows(all)

```

In []:

```

# Downloading API urls in json format to the Local DB

import requests
import os
import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-105hhrg40050']
count = 0
for jacket in df1['Filename']:

    try:
        #print (set(os.listdir(APIs)))
        #print (jacket+".json")
        if jacket+".json" not in set(os.listdir(APIs)):
            url = 'https://api.govinfo.gov/packages/'+jacket+'/mods?&api_key=XNEgGxjbEs
zIMyIeni9xpgdkqy60QD5p9S4Vvdlc'

            r = requests.get(url)

            with open('data.xml', 'w' , encoding="utf8") as f:
                f.write(r.text)

            with open("data.xml", 'r' , encoding="utf8") as f:
                xmlString = f.read()

            #print ("XML input (data.xml):")
            #print(xmlString)

            jsonString = json.dumps(xmltodict.parse(xmlString), indent=4)

            jsonObj = json.loads(jsonString)

            #print("\nJSON output(output.json):")
            #print(jsonString)

            file = APIs + jacket+ ".json"

            with open(file, 'w', encoding="utf8") as f:
                f.write(jsonString)

```

```
except:
    print(jacket)
```

In []:

```
# Downloading full text in .txt format to the local DB

import os
import urllib.request
import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

FullText = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/FullTexts/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-105hhrg40050']
count = 0
for jacket in df1['Filename']:

    try:
        #print (set(os.listdir(APIs)))
        #print (jacket+".json")
        if jacket+".txt" not in set(os.listdir(FullText)):

            url = 'https://api.govinfo.gov/packages/'+jacket+'/granules/'+jacket+'/htm?api_key=XNEgGxjbEszIMyIeni9xpgdkqy60QD5p9S4Vvd1c'

            file = FullText + jacket + ".txt"

            urllib.request.urlretrieve(url, file)

    except:
        print(jacket)
```

In []:

```
# Read the file in local DB

file = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/FullTexts/CHRG-115hhrg23826.txt"

file_lines = open(file).readlines()
print (file_lines[:20])
```


In []:

```
# Individual CSVs
# heldDate extraction

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

for file in os.listdir(results_csvs):

    try:

        #with open(APIs+file, 'r') as f:
        #    xmlString = f.read()

        #print ("XML input (data.xml):")
        #print(xmlString)

        file = file.replace('.csv', '.json')

        with open(APIs+file) as data_file:
            jsonObj = json.load(data_file)
            #print(jsonObj)

        file = file.replace('.json', '.csv')

        # if file == 'CHRG-100shrg83712.csv' or file == 'CHRG-102hhrg67539.csv' or file
        == 'CHRG-103hhrg66111.csv' :
            # continue
            with open(results_csvs+file, 'r', encoding="utf8") as csvinput:
                with open(results_csvs_new+file, 'w+', encoding="utf8") as csvoutput:
                    writer = csv.writer(csvoutput, lineterminator='\n')
                    reader = csv.reader(csvinput)

                    all = []
                    row = next(reader)

                    row.append('heldDate')
                    all.append(row)
                    #print (row)
                    #try:

                    for row in reader:
                        try:
                            heldDate = []
                            added = False
                            exists = False

                            for item in (jsonObj["mods"]["extension"]):
                                #for item in extension:
                                #print (item)
                                if "heldDate" in item:
                                    exists = True
                                    if isinstance(item["heldDate"], list):
                                        for date in item["heldDate"]:
                                            heldDate.append(date)
                                            added = True
                                            #print (heldDate)
                                else:
                                    row.append(item["heldDate"])

```

```
                                #print (item["heldDate"])
                                break

        if exists == False:
            row.append("-")
        if added:
            row.append("; \n".join(heldDate))
            #break

    except:
        row.append("-")

    all.append(row)

#except:
#    writer.writerow(all)
#    continue
writer.writerow(all)

except:
    print (file)
```

In []:

```
# metadata_results
# heIdDate extraction

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        try:

            all = []
            row = next(reader)

            row.append('heldDate')
            all.append(row)
            #print (row)
            #try:

            for row in reader:
                try:

                    #if (row[6].strip()+'.csv' in os.listdir(results_csv
s)):

                        #print (row[6].strip()+'.csv')
                        file = row[6].strip()

                        file = file + '.json'

                        with open(APIs+file) as data_file:
                            jsonObj = json.load(data_file)
                            #print(jsonObj)

                            heldDate = []
                            added = False
                            exists = False

                            for item in (jsonObj["mods"]["extension"]):
                                #for item in extension:
                                    #print (item)
                                    if "heldDate" in item:
                                        exists = True
                                        if isinstance(item["heldDate"], lis
t):

                                            for date in item["heldDate"]:
                                                heldDate.append(date)
                                                added = True
                                                #print (heldDate)
                                            else:
                                                row.append(item["heldDate"])
                                                # print (item["heldDate"])
                                                break

```

```
        if exists == False:
            row.append("-")
        if added:
            row.append("; \n".join(heldDate))
            #break

    except:
        row.append("-")

    all.append(row)

#except:
#    writer.writerow(all)
#    continue
writer.writerow(all)

except:
    print (file)
```

In []:

```
# GPO agencies for sample 500 CSVs
# Individual CSVs

import csv
import pandas as pd
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

df1 = pd.read_csv(sample500)
#print(df1['filename'])

df = pd.read_csv(gpo2)
agencies = []
for i in df['Agency']:
    temp = i.replace('U.S.', 'United States')
    temp = temp.replace('U.S', 'United States')
    temp = temp.replace('Dep.', 'Department')

    agencies.append(temp)

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')

for file in df1['filename']:
    try:

        #print ( row1['Affiliation'] + ' : ' + agency + '\t' + str(max_score))

        with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
            with open(sample500GPOOutput+file+'.csv', 'w+', encoding="utf8") as
```


csvoutput:

```

writer = csv.writer(csvoutput, lineterminator='\n')
reader = csv.reader(csvinput)

all = []
row = next(reader)

row.append('Government agencies')
all.append(row)
#print (row)
#try:

for row in reader:

    file1 = pd.read_csv(results_csvs + file + '.csv')

    max_score = 0
    agency = '-'
    #print (row[18])
    if (str(row[16]).strip() == "Yes"):
        max_score = 0
        agency = '-'
        for i in (set(agencies)):
            score = fuzz.token_set_ratio( i.lower(), row[18
].lower())

            if (score > max_score):
                max_score = score
                agency = i

    if max_score == 100:
        row.append(agency)
    else:
        row.append(agency)

    all.append(row)

#except:
#    writer.writerow(all)
#    continue
writer.writerow(all)

except:
    print (file)

```

In []:

```
# Sentiment analysis for sample 500 CSVs
# Individual CSVs

import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer
nltk.download('vader_lexicon')

import csv
import pandas as pd
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

sample500SAOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-SA/"

df1 = pd.read_csv(sample500)
#print(df1['filename'])

sid = SentimentIntensityAnalyzer()

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')

for file in df1['filename']:

    try:

        #print ( row1['Affiliation'] + ' : ' + agency + '\t' + str(max_score))

        with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
            with open(sample500SAOutput+file+'.csv', 'w+', encoding="utf8") as csvoutput:
```

```
writer = csv.writer(csvoutput, lineterminator='\n')
reader = csv.reader(csvinput)

all = []
row = next(reader)

row.append('Sentiment analysis')
all.append(row)
#print (row)
#try:

#print (ss)
#print (max(ss, key=ss.get))
#break
for row in reader:

    #df2 = pd.read_csv(results_csvs+file+'.csv')

    #print (df2['cleaned'])
    ss = sid.polarity_scores(row[12])

    del (ss['compound'])

    #print (row[12])

    if ( max(ss, key=ss.get) == 'neu'):
        row.append('Neutral')

    if ( max(ss, key=ss.get) == 'neg'):
        row.append('Negative')

    if ( max(ss, key=ss.get) == 'pos'):
        row.append('Positive')

    all.append(row)

#except:
#    writer.writerow(all)
#    continue
writer.writerow(all)

except:
    print (file)
```

In []:

```

# Metadata_results
# Witness names & Affiliations, Members of the congress from FULL Texts - Scrapped Wit
nesses

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hear
ings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional
_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hear
ings_from_server/gpo_tools/results_csvs_new/"

FullTexts = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_f
rom_server/gpo_tools/FullTexts/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhr27211']
count = 0

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

```

```

countWitness = 0
with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Scrapped witnesses')
        all.append(row)

        for row in reader:

            #try:
            if (row[9]!='Appropriation' and row[9]!='Nomination') and row[14]!='-' and row[16]!='Yes' and row[14]!='-':

                #print(row)
                if row[6]+'.txt' in set(os.listdir(FullTexts)):
                    filename = FullTexts+row[6]+'.txt'
                    lines = open(filename, "r", encoding="utf8").readlines()

                    #print (lines)
                    strippedLines = []
                    for line in lines:
                        #print (line.strip())
                        strippedLines.append(line.strip())

                    if ('CONTENTS' in strippedLines and 'Statement of:' in strippedLines):

                        startingIndex = strippedLines.index('Statement of:')

                        #print (startingIndex)
                        witness = []

                        #print ('\n'+row[6])
                        #print (lines)
                        witnessStr = []
                        firstHit = 0
                        for i in range(startingIndex+1, len(strippedLines)):

                            if ' ' in lines[i]:

                                if lines[i].strip() == '':
                                    continue

                                if re.search(r"\.(\.)+( *)[0-9]", lines[i]):

                                    if (firstHit == 0):
                                        x = re.sub('\.(\.)+( *)', '', lines[i])

                                        witness.append(x.strip())

                                        firstHit = 1

                                    elif ';' in lines[i]:

```

```

ip().split(';')[0]+'\\n')
ip().split(';')[1].replace('and', '').strip()

ip()+ ' ')

else:
    witness.append(lines[i].str

firstHit = 0

else:
    break

#print ("".join(witness))
if(len("".join(witness)) < 6000):
    row.append("".join(witness))
    row[14] = 'Refer column S'
    countWitness += 1

elif ('C O N T E N T S' in strippedLines and 'S
TATEMENTS' in strippedLines):

    #countWitness += 1

    startingIndex = strippedLines.index('ST
ATEMENTS')

    #print (startingIndex)
    witness = []

    #print ('\\n'+row[6])
    #print (lines)
    witnessStr = []
    firstHit = 0
    for i in range(startingIndex+1, len(lin
es)):

        if 'APPENDIX' in lines[i] or 'A
ppendix' in lines[i]:

            break
        if 'Page' in lines[i]:
            continue

        if lines[i].isupper():
            break

        if lines[i].strip() == '':
            continue

        if re.search(r"\\.(\\.)+( *)[0-9]
*(\\*)*$",lines[i]):

            if(firstHit == 0):
                x = re.sub('\\.(\\.)+( *)
[0-9]*(\\*)*$', '', lines[i])

                witness.append(x.strip

                firstHit = 1

            elif ';' in lines[i]:
                witness.append(lines[i].spl
it(';')[0].strip()+ '\\n')

```

```

it(';')[1].replace('and','').strip())

ip()+' ')

Y' in strippedLines):

STIMONY')

es)):

ppendix' in lines[i]:

*(\*)*$",lines[i]):

[0-9]*(\*)*$', '', lines[i])

()+'\n')

it(';')[0].strip()+'\n')

it(';')[1].replace('and','').strip())

witness.append(lines[i].spl

firstHit = 0

else:
witness.append(lines[i].str

firstHit = 0

#print ("".join(witness))
if(len("".join(witness)) < 6000):
    row.append("".join(witness))
    row[14] = 'Refer column S'
    countWitness += 1

elif ('CONTENTS' in strippedLines and 'TESTIMON

#countWitness += 1

startingIndex = strippedLines.index('TE

#print (startingIndex)
witness = []

#print ('\n'+row[6])
#print (lines)
witnessStr = []
firstHit = 0
for i in range(startingIndex+1, len(lin

if 'APPENDIX' in lines[i] or 'A

break
if 'Page' in lines[i]:
    continue

if lines[i].isupper():
    break

if lines[i].strip() == '':
    continue

if re.search(r"\.(\\.)+( *)[0-9]

if(firstHit == 0):
    x = re.sub('\\.(\\.)+( *)

witness.append(x.strip

firstHit = 1

elif ';' in lines[i]:
    witness.append(lines[i].spl

witness.append(lines[i].spl

firstHit = 0

else:

```

```

ip()+' ')

witness.append(lines[i].str

firstHit = 0

#print ("".join(witness))
if(len("".join(witness)) < 6000):
    row.append("".join(witness))
    row[14] = 'Refer column S'
    countWitness += 1

elif ('C O N T E N T S' in strippedLines and 'T
estimony of:' in strippedLines):
    #countWitness += 1

    startingIndex = strippedLines.index('Te
stimony of:')

#print (startingIndex)
witness = []

#print ('\n'+row[6])
#print (lines)
witnessStr = []
firstHit = 0
for i in range(startingIndex+1, len(lin
es)):

    if 'APPENDIX' in lines[i] or 'A
ppendix' in lines[i]:
        break
    if 'Page' in lines[i]:
        continue

    if lines[i].isupper():
        break

    if lines[i].strip() == '':
        continue

    if re.search(r"\.(\.)+( *)[0-9]
*(\*)*$",lines[i]):
        if(firstHit == 0):
            x = re.sub('\.(\.)+( *)
[0-9]*(\*)*$','', lines[i])

            witness.append(x.strip

            firstHit = 1

        elif ';' in lines[i]:
            witness.append(lines[i].spl

            witness.append(lines[i].spl

            firstHit = 0

        else:
            witness.append(lines[i].str

            firstHit = 0

ip()+' ')

```



```

        #print ("".join(witness))
        if(len("".join(witness)) < 6000):
            row.append("".join(witness))
            row[14] = 'Refer column S'
            countWitness += 1

    elif ('C O N T E N T S' in strippedLines and 'C
HRONOLOGICAL LIST OF WITNESSES' in strippedLines):
        #countWitness += 1

        startingIndex = strippedLines.index('CH
RONOLOGICAL LIST OF WITNESSES')

        #print (startingIndex)
        witness = []

        #print ('\n'+row[6])
        #print (lines)
        witnessStr = []
        firstHit = 0
        for i in range(startingIndex+1, len(lin
es)):

            if lines[i].isupper():
                break

            if lines[i].strip() == '':
                continue

            if re.search(r"\.(\\.)+( *)[0-9]
*(\\)**$",lines[i]):

                if(firstHit == 0):
                    x = re.sub('\\.(\\.)+( *)
[0-9]*(\\)*$','', lines[i])

                    witness.append(x.strip
                    ())+'\\n')

                    firstHit = 1

            elif ';' in lines[i]:
                witness.append(lines[i].spl

                witness.append(lines[i].spl

                firstHit = 0

            else:
                witness.append(lines[i].str

                firstHit = 0

        #print ("".join(witness))
        if(len("".join(witness)) < 6000):
            row.append("".join(witness))
            row[14] = 'Refer column S'
            countWitness += 1

    elif ('C O N T E N T S' in strippedLines and (
'Panel I' in strippedLines or 'PANEL I' in strippedLines)) :
        #countWitness += 1

```

```
('Panel I')
```

```
('PANEL I')
```

```
es)):
```

```
-':
```

```
ppendix' in lines[i]:
```

```
*(\*)*$",lines[i]):
```

```
[0-9]*(\*)*$',' ', lines[i])
```

```
()+'\n')
```

```
it(';')[0].strip()+'\n')
```

```
it(';')[1].replace('and','').strip())
```

```
ip()+' ')
```

```
if 'Panel I' in strippedLines:
    startingIndex = strippedLines.index
```

```
if 'PANEL I' in strippedLines:
    startingIndex = strippedLines.index
```

```
#print (startingIndex)
```

```
witness = []
```

```
#print ('\n'+row[6])
```

```
#print (lines)
```

```
witnessStr = []
```

```
firstHit = 0
```

```
for i in range(startingIndex+1, len(lin
```

```
if lines[i].strip == '-----
```

```
break
```

```
if lines[i].strip() == '':
```

```
continue
```

```
if 'Panel' in lines[i]:
```

```
continue
```

```
if 'APPENDIX' in lines[i] or 'A
```

```
break
```

```
if 'Page' in lines[i]:
```

```
continue
```

```
if lines[i].isupper():
```

```
break
```

```
if re.search(r"\.(\.)+( *)[0-9]
```

```
if(firstHit == 0):
```

```
x = re.sub('\.(\.)+( *)
```

```
witness.append(x.strip
```

```
firstHit = 1
```

```
elif ';' in lines[i]:
```

```
witness.append(lines[i].spl
```

```
witness.append(lines[i].spl
```

```
firstHit = 0
```

```
else:
```

```
witness.append(lines[i].str
```

```
firstHit = 0
```

```
#print ("".join(witness))
```

```
if(len("".join(witness)) < 6000):
```

```
row.append("".join(witness))
```



```

#print ("".join(witness))
if(len("".join(witness)) < 6000):
    row.append("".join(witness))
    row[14] = 'Refer column S'
    countWitness += 1

elif ('C O N T E N T S' in strippedLines and 'S
tatements:' in strippedLines):

    #countWitness += 1

    startingIndex = strippedLines.index('St
atements:')

    #print (startingIndex)
    witness = []

    #print ('\n'+row[6])
    #print (lines)
    witnessStr = []
    firstHit = 0
    for i in range(startingIndex+1, len(lin
es)):

        if 'APPENDIX' in lines[i] or 'A
ppendix' in lines[i]:

            break
        if 'Page' in lines[i]:
            continue

        if lines[i].isupper():
            break

        if lines[i].strip() == '':
            continue

        if re.search(r"\.(\.)+( *)[0-9]
*(\*)*$",lines[i]):

            if(firstHit == 0):
                x = re.sub('\.(\.)+( *)
[0-9]*(\*)*$', '', lines[i])

                witness.append(x.strip
                ())+'\n')

                firstHit = 1

            elif ';' in lines[i]:
                witness.append(lines[i].spl
                it(';')[0].strip()+'\n')

                witness.append(lines[i].spl
                it(';')[1].replace('and','').strip())

                firstHit = 0

            else:
                witness.append(lines[i].str
                ip()+' ')

                firstHit = 0

    #print ("".join(witness))
    if(len("".join(witness)) < 6000):
        row.append("".join(witness))

```

```

row[14] = 'Refer column S'
countWitness += 1

elif ('C O N T E N T S' in strippedLines and 'W
ITNESS' in strippedLines):

    #countWitness += 1

    startingIndex = strippedLines.index('WI
TNESS')

    #print (startingIndex)
    witness = []

    #print ('\n'+row[6])
    #print (lines)
    witnessStr = []
    firstHit = 0
    for i in range(startingIndex+1, len(lin
es)):

        if 'APPENDIX' in lines[i] or 'A
ppendix' in lines[i]:

            break
        if 'Page' in lines[i]:
            continue

        if lines[i].isupper():
            break

        if lines[i].strip() == '':
            continue

        if re.search(r"\.(\.)+( *)[0-9]
*(\*)*$",lines[i]):

            if(firstHit == 0):
                x = re.sub('\.(\.)+( *)
[0-9]*(\*)*$', '', lines[i])

                witness.append(x.strip
                ())+'\n')

                firstHit = 1

            elif ';' in lines[i]:
                witness.append(lines[i].spl
                witness.append(lines[i].spl
                firstHit = 0

            else:
                witness.append(lines[i].str
                firstHit = 0

    #print ("".join(witness))
    if(len("".join(witness)) < 6000):
        row.append("".join(witness))
        row[14] = 'Refer column S'
        countWitness += 1

elif ('C O N T E N T S' in strippedLines and 'W
itnesses:' in strippedLines):

```

tnesses:')

es)):

ppendix' in lines[i]:

()*\$",lines[i]):

[0-9]*(*)*\$', '', lines[i])

()+'\n')

it(';')[0].strip()+'\n')

it(';')[1].replace('and','').strip())

ip()+' ')

TNERS FOR COOPERATION' in strippedLines and 'WITNESSES' in strippedLines):

TNESSES')

#countWitness += 1

startingIndex = strippedLines.index('Wi

#print (startingIndex)

witness = []

#print ('\n'+row[6])

#print (lines)

witnessStr = []

firstHit = 0

for i in range(startingIndex+1, len(lin

if 'APPENDIX' in lines[i] or 'A

break

if 'Page' in lines[i]:

continue

if lines[i].isupper():

break

if lines[i].strip() == '':

continue

if re.search(r"\.(\\.)+(*)[0-9]

if(firstHit == 0):

x = re.sub('\\.(\\.)+(*)[0-9]

witness.append(x.strip

firstHit = 1

elif ';' in lines[i]:

witness.append(lines[i].spl

witness.append(lines[i].spl

firstHit = 0

else:

witness.append(lines[i].str

firstHit = 0

#print ("".join(witness))

if(len("".join(witness)) < 6000):

row.append("".join(witness))

row[14] = 'Refer column S'

countWitness += 1

elif ('THE FUTURE OF THE OSCE MEDITERRANEAN PAR

#countWitness += 1

startingIndex = strippedLines.index('WI

#print (startingIndex)

```
es)):
```

```
ppendix' in lines[i]:
```

```
*(\*)*$",lines[i]):
```

```
[0-9]*(\*)*$', '', lines[i])
```

```
()+'\n')
```

```
it(';')[0].strip()+'\n')
```

```
it(';')[1].replace('and','').strip())
```

```
ip()+' ')
```

```
age' in strippedLines):
```

```
ge')
```

```
witness = []
```

```
#print ('\n'+row[6])
```

```
#print (lines)
```

```
witnessStr = []
```

```
firstHit = 0
```

```
for i in range(startingIndex+1, len(lin
```

```
if 'APPENDIX' in lines[i] or 'A
```

```
break
```

```
if 'Page' in lines[i]:
```

```
continue
```

```
if lines[i].isupper():
```

```
break
```

```
if lines[i].strip() == '':
```

```
continue
```

```
if re.search(r"\.(\\.)+( *)[0-9]
```

```
if(firstHit == 0):
```

```
x = re.sub('\\.(\\.)+( *)
```

```
witness.append(x.strip
```

```
firstHit = 1
```

```
elif ';' in lines[i]:
```

```
witness.append(lines[i].spl
```

```
witness.append(lines[i].spl
```

```
firstHit = 0
```

```
else:
```

```
witness.append(lines[i].str
```

```
firstHit = 0
```

```
#print ("".join(witness))
```

```
if(len("".join(witness)) < 6000):
```

```
row.append("".join(witness))
```

```
row[14] = 'Refer column S'
```

```
countWitness += 1
```

```
elif ('C O N T E N T S' in strippedLines and 'P
```

```
#countWitness += 1
```

```
startingIndex = strippedLines.index('Pa
```

```
#print (startingIndex)
```

```
witness = []
```

```
#print ('\n'+row[6])
```

```
#print (lines)
```

```
witnessStr = []
```

```

es)):
    firstHit = 0
    for i in range(startingIndex+1, len(lines)):
        if 'APPENDIX' in lines[i] or 'A' in lines[i]:
            break
        if 'Page' in lines[i]:
            continue
        if lines[i].isupper():
            break
        if lines[i].strip() == '':
            continue
        if re.search(r"\.(\\.)+( *)[0-9]", lines[i]):
            if firstHit == 0:
                x = re.sub('\\\\.(\\.)+( *)', '\\.\\.', lines[i])
                witness.append(x.strip())
                firstHit = 1
            elif ';' in lines[i]:
                witness.append(lines[i].split(';')[0].strip())
                witness.append(lines[i].split(';')[1].replace('and', '').strip())
                firstHit = 0
            else:
                witness.append(lines[i].strip())
                firstHit = 0

    #print ("".join(witness))

    if(len("".join(witness)) < 6000):
        row.append("".join(witness))
        row[14] = 'Refer column S'
        countWitness += 1
    else:
        row.append('-')
        #row[14] = 'Refer column S'

    else:
        row.append('-')

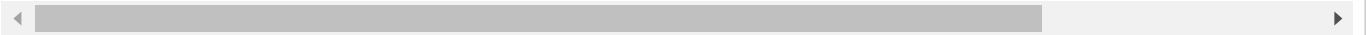
    #if countWitness !=0:
    #    break
    #except:
    #    row.append("-")
    #    row.append("-")
    all.append(row)
    #    all.append(row)
    #    continue
except:

```



```
#     writer.writerows(all)
#     continue
writer.writerows(all)

print(countWitness)
```



In []:

```
# Metadata_results
# Witness names & Affiliations, Members of the congress from FULL Texts - Scrapped Withn
esses for individual CSVs

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

FullTexts = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/FullTexts/"

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hrg27211']
count = 0

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))
```

```

countWitness = 0

gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

sample500SAOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-SA/"

#df1 = pd.read_csv(sample_csvs_new)
#print(df1['filename'])

#sid = SentimentIntensityAnalyzer()

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')

scrappedWD = {}

with open(metadata_results, 'r', encoding="utf8") as csvinput:

    reader = csv.reader(csvinput)

    for row in reader:

        #try:
        if row[13] == 'Refer column R' :

            scrappedWD[row[6]] = row[17]

for k, v in scrappedWD.items():

```

```
print (v.split('\n'))  
break
```

In []:

```
# Metadata_results
# Witness names & Affiliations, Members of the congress from FULL Texts - Scrapped Withn
esses for individual CSVs

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

FullTexts = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/FullTexts/"

#df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhr27211']
count = 0

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))
```

```
countWitness = 0
```

```
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
```

```
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"
```

```
metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
```

```
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"
```

```
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
```

```
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"
```

```
sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
```

```
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"
```

```
sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"
```

```
sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"
```

```
sample500SAOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-SA/"
```

```
df1 = pd.read_csv(sample500)
```

```
#print(df1['filename'])
```

```
#sid = SentimentIntensityAnalyzer()
```

```
#print (set(agencies))
```

```
from fuzzywuzzy import fuzz
```

```
from fuzzywuzzy import process
```

```
#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')
```

```
for file in set(os.listdir(results_csvs)):
```

```
#for file in set(os.listdir(results_csvs)):
```

```
    #print (file)
```

```
    #print (set(os.listdir(results_csvs)))
```

```
    #file = file + '.csv'
```

```
    #if file in set(os.listdir(results_csvs)):
```

```
        file = file.replace('.csv', '')
```

```
        with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
```

```
            with open(results_csvs_new+file+'.csv', 'w+', encoding="utf8") as csvoutput
```

```
            :
```

```
                writer = csv.writer(csvoutput, lineterminator='\n')
```

```
                reader = csv.reader(csvinput)
```

```
                all = []
```

```

row = next(reader)

row.append('Scrapped witnesses')
all.append(row)

for row in reader:
    hit = 0
    if row[7].strip() in scrappedWD.keys():
        tempWit = scrappedWD[row[7].strip()]
        #print (tempWit)
        name = row[3] + ' ' + row[5] + ' ' + row[17]
        for j in tempWit.split('\n'):
            if fuzz.token_sort_ratio("".join(j.lower().split())[:4
]), name.lower()) > 40 and j.strip()!='':
                row.append(j.strip())
                row[16] = 'Yes'
                hit = 1
                #break
                #print (fuzz.token_sort_ratio("".join(j.lower().spl
it())[:4]), name.lower()))

                #print ("".join(j.lower().split())[:4]))
                #print (name.lower())
                break

        if hit == 0:
            row.append('-')
    else:
        row.append('-')

    all.append(row)

writer.writerows(all)

```

In []:

```

# Cleaning witness, scrapped witness column

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

for file in set(os.listdir(results_csvs)):
    with open(results_csvs+file, 'r', encoding="utf8") as csvinput:
        with open(results_csvs_new+file, 'w+', encoding="utf8") as csvoutput:
            writer = csv.writer(csvoutput, lineterminator='\n')
            reader = csv.reader(csvinput)

            all = []
            row = next(reader)

            all.append(row)

            for row in reader:

                if row[18].strip() == 'United States Senate' or row[18].strip() == 'United States House of Representatives':
                    row[16] = 'No'
                    row[20] = '-'

                all.append(row)

```



```
writer.writerow(all)

print ('asdf')
```

In []:

```
# Creating dictionary of acronyms and agencies

import os

import math
import csv
import pandas as pd

gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

#print(df1['filename'])

df = pd.read_csv(gpo2)
agencies = []
acronyms = []

acroMap = {}

for i in (df['Agency']):
    agencies.append(i)

for i in (df['Alternate Name']):
    acronyms.append(i)

for i in acronyms:
    if not(pd.isnull(i)):
        index = acronyms.index(i)
        acroMap[i] = agencies[index]

print((acroMap.keys()))
```

In []:

```
# Creating dictionary of acronyms and states

import os

import math
import csv
import pandas as pd

usstates = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/us_states.csv"

df = pd.read_csv(usstates, header=None)
states= []
acronyms = []

acroMapStates = {}

for i in (df.iloc[:,1]):
    states.append(i)

for i in (df.iloc[:,2]):
    acronyms.append(i)

for i in acronyms:
    #if not(pd.isnull(i)):
        index = acronyms.index(i)
        acroMapStates[i] = states[index]

print((acroMapStates.keys()))
```

In []:

```
# GPO agencies for individual CSVs
# Exact matching on agency names and acronyms, states, Inspector General

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math
import csv
import pandas as pd
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"
results_csvs_new1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new1/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

df1 = pd.read_csv(sample500)
#print(df1['filename'])
```

```

df = pd.read_csv(gpo2)
agencies = []
for i in (df['Agency']):
    temp = i.replace('U.S.', 'United States')
    temp = temp.replace('US', 'United States')
    temp = temp.replace('Dep.', 'Department')
    temp = temp.replace('Dept.', 'Department')
    temp = temp.replace('Dept', 'Department')
    temp = temp.replace('Assoc', 'Association')
    temp = temp.replace('Assoc.', 'Association')
    temp = temp.replace('Brd', 'Board')
    temp = temp.replace('Brd.', 'Board')
    temp = temp.replace('DC', 'District of Columbia')
    temp = temp.replace('D.C.', 'District of Columbia')

    temp = temp.replace('.', ' ')
    temp = temp.replace('; ', ' ')
    temp = temp.replace('-', ' ')
    temp = temp.replace(':', ' ')
    temp = temp.replace('.', ' ')

    temp = temp.replace('.', '')

    for i in temp.split():
        if i in acroMap.keys():
            temp = temp.replace(i, acroMap[i])

    for i in temp.split():
        if i in acroMapStates.keys():
            temp = temp.replace(i, acroMapStates[i])

    agencies.append(temp)

JK = []
UA = []
Parent = []

for i in (df['JK Code']):
    JK.append(i)
for i in (df['UA Code']):
    UA.append(i)
for i in (df['Parent UA Code']):
    Parent.append(i)

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')

#for file in df1['filename']:

#    try:

        #print ( row1['Affiliation'] + ' : ' + agency + '\t' + str(max_score))
#agencies = agencies[:100]
for file in set(os.listdir(results_csvs)):
    if file not in set(os.listdir(results_csvs_new)):

```

```

with open(results_csvs+file,'r', encoding="utf8") as csvinput:
    with open(results_csvs_new+file, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

    all = []
    row = next(reader)

    row.append('Agency')

    row.append('JK code')
    row.append('UA code')
    row.append('Parent UA code')
    row.append('US State')
    row.append('Inspector General')
    all.append(row)
    #print (row)
    #try:

    for row in reader:

        if row[16] == 'Yes':
            max_score = 0
            agency = '-'
            jk = '-'
            ua = '-'
            parent = '-'
            aff = row[18] + ' '+row[20]

            aff = aff.replace('U.S.', 'United States')
            aff = aff.replace('US', 'United States')
            aff = aff.replace('Dep.', 'Department')
            aff = aff.replace('Dept.', 'Department')
            aff = aff.replace('Dept', 'Department')
            aff = aff.replace('Assoc', 'Association')
            aff = aff.replace('Assoc.', 'Association')
            aff = aff.replace('Brd', 'Board')
            aff = aff.replace('Brd.', 'Board')
            aff = aff.replace('DC', 'District of Columbia')
            aff = aff.replace('D.C.', 'District of Columbia')

            aff = aff.replace('.',',')
            aff = aff.replace(';',' ')
            aff = aff.replace('-',',')
            aff = aff.replace(':',',')
            aff = aff.replace('..',',')

            aff = aff.replace('.',',')

            for i in aff.split():
                if i in acroMap.keys():
                    aff = aff.replace(i,acroMap[i])

            for i in aff.split():
                if i in acroMapStates.keys():
                    aff = aff.replace(i,acroMapStates[i])

            hit = 0
            for i in ((agencies)):
                #score = fuzz.WRatio( i, aff )

```

```

    #if (score > max_score):
    if i in aff:
        #max_score = score
        agency = i
        index = agencies.index(i)
        jk = JK[index]
        ua = UA[index]
        parent = Parent[index]

        row.append(agency)
        row.append(jk)
        row.append(ua)
        row.append(parent)
        hit = 1
        break
    ...
    if max_score >= 90:
        row.append(agency)
        row.append(jk)
        row.append(ua)
        row.append(parent)

    else:
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')
    ...

    if hit == 0:
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')

    states = 0

    for i in acroMapStates.values():
        if i in aff:
            row.append(i)
            states = 1
            break

    if states == 0:
        row.append('-')

    if 'IG' in aff or 'Inspector General' in aff or 'Inspe
c. General' in aff:
        row.append('Yes')
    else:
        row.append('No')

    else:
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')

    all.append(row)

```

```
#except:  
#     writer.writerow(all)  
#     continue  
writer.writerow(all)
```

In []:

```
from fuzzywuzzy import fuzz  
from fuzzywuzzy import process  
print (fuzz.partial_ratio( 'Hon. Peter J. Visclosky, a Representative in Congress from  
the State of Indiana', 'v'))
```

In []:

```
# GPO agencies for metadata
# Exact matching on agency names and acronyms, states, Inspector General

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math
import csv
import pandas as pd
gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"
results_csvs_new1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new1/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

df1 = pd.read_csv(sample500)
#print(df1['filename'])
```



```

df = pd.read_csv(gpo2)
agencies = []
for i in (df['Agency']):
    temp = i.replace('U.S.', 'United States')
    temp = temp.replace('US', 'United States')
    temp = temp.replace('Dep.', 'Department')
    temp = temp.replace('Dept.', 'Department')
    temp = temp.replace('Dept', 'Department')
    temp = temp.replace('Assoc', 'Association')
    temp = temp.replace('Assoc.', 'Association')
    temp = temp.replace('Brd', 'Board')
    temp = temp.replace('Brd.', 'Board')
    temp = temp.replace('DC', 'District of Columbia')
    temp = temp.replace('D.C.', 'District of Columbia')

    temp = temp.replace('.', ' ')
    temp = temp.replace('; ', ' ')
    temp = temp.replace('-', ' ')
    temp = temp.replace(':', ' ')
    temp = temp.replace('.', ' ')

    temp = temp.replace('.', '')

    for i in temp.split():
        if i in acroMap.keys():
            temp = temp.replace(i, acroMap[i])

    for i in temp.split():
        if i in acroMapStates.keys():
            temp = temp.replace(i, acroMapStates[i])

    agencies.append(temp)

JK = []
UA = []
Parent = []

for i in (df['JK Code']):
    JK.append(i)
for i in (df['UA Code']):
    UA.append(i)
for i in (df['Parent UA Code']):
    Parent.append(i)

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')

#for file in df1['filename']:

#    try:

        #print ( row1['Affiliation'] + ' : ' + agency + '\t' + str(max_score))
#agencies = agencies[:100]
#for file in set(os.listdir(results_csvs)):
#    if file not in set(os.listdir(results_csvs_new)):

```

```

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Agency')

        row.append('JK code')
        row.append('UA code')
        row.append('Parent UA code')
        row.append('US State')
        row.append('Inspector General')
        all.append(row)
        #print (row)
        #try:

        for row in reader:

            if row[13].strip() != '-':
                max_score = 0
                agency = '-'
                jk = '-'
                ua = '-'
                parent = '-'

                agencyL = []
                jkL = []
                uaL = []
                parentL = []
                stateL = []
                IGL = []

                if row[13] == 'Refer column R':
                    affs = row[17].split('\n')
                else:
                    affs = row[13].split('\n')

                for aff in affs:
                    if aff.strip() != '':
                        aff = aff.replace('U.S.', 'United States')
                        aff = aff.replace('US', 'United States')
                        aff = aff.replace('Dep.', 'Department')
                        aff = aff.replace('Dept.', 'Department')
                        aff = aff.replace('Dept', 'Department')
                        aff = aff.replace('Assoc', 'Association')
                        aff = aff.replace('Assoc.', 'Association')
                        aff = aff.replace('Brd', 'Board')
                        aff = aff.replace('Brd.', 'Board')
                        aff = aff.replace('DC', 'District of Columbia')
                        aff = aff.replace('D.C.', 'District of Columbi

a')

                        aff = aff.replace('.', ',')
                        aff = aff.replace(';', ',')
                        aff = aff.replace('-', ',')
                        aff = aff.replace(':', ',')
                        aff = aff.replace('.', ',')

```

```

aff = aff.replace('.', '')

for i in aff.split():
    if i in acroMap.keys():
        aff = aff.replace(i, acroMap[i])

for i in aff.split():
    if i in acroMapStates.keys():
        aff = aff.replace(i, acroMapStates[i])

hit = 0
for i in ((agencies)):
    #score = fuzz.WRatio( i, aff )
    #if (score > max_score):
    if i in aff:
        #max_score = score
        agency = i
        index = agencies.index(i)
        jk = JK[index]
        ua = UA[index]
        parent = Parent[index]

        agencyL.append(str(agency))
        jkL.append(str(jk))
        uaL.append(str(ua))
        parentL.append(str(parent))
        hit = 1
        break
    ...
    if max_score >= 90:
        row.append(agency)
        row.append(jk)
        row.append(ua)
        row.append(parent)

    else:
        row.append('-')
        row.append('-')
        row.append('-')
        row.append('-')
    ...

if hit == 0:
    agencyL.append('-')
    jkL.append('-')
    uaL.append('-')
    parentL.append('-')

states = 0

for i in acroMapStates.values():
    if i in aff:
        stateL.append(i)
        states = 1
        break

if states == 0:
    stateL.append('-')

if 'IG' in aff or 'Inspector General' in aff or

```

```
'Inspec. General' in aff:

    IGL.append('Yes')
else:
    IGL.append('No')

    row.append("\n".join(agencyL))
    row.append("\n".join(jkL))
    row.append("\n".join(uaL))
    row.append("\n".join(parentL))
    row.append("\n".join(stateL))
    row.append("\n".join(IGL))

else:
    row.append('- ')
    row.append('- ')
    row.append('- ')
    row.append('- ')
    row.append('- ')
    row.append('- ')

    all.append(row)

except:
#     writer.writerow(all)
#     continue
writer.writerow(all)
```

In []:

```

# Adding "Bills" column in all individual CSVs

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

FullTexts = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/FullTexts/"

#df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

#files = set(os.listdir(results_csvs)) - set(os.listdir(results_csvs_new))

countWitness = 0

```

```

gpo = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019.csv"
gpo2 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Master agencies list_Feb. 2019_v2.csv"

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

sample500 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/Sample_500_108th-112th_Congresses_1.31.19.csv"

sample500GPOOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-GPOs/"

sample500SAOutput = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs500-SA/"

sampleBill = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/Extras/CHRG-109shrg26254_coded_bills (1).csv"

#df1 = pd.read_csv(sample500)
#print(df1['filename'])

#sid = SentimentIntensityAnalyzer()

#print (set(agencies))

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#file = pd.read_csv(sample_csvs + 'CHRG-105hhrg40051' + '.csv')
count = 0
#for file in set(os.listdir(results_csvs)):
for file in set(os.listdir(results_csvs)):
    #print (file)
    #print (set(os.listdir(results_csvs)))
    #file = file + '.csv'
    #if file in set(os.listdir(results_csvs)):
    file = file.replace('.csv', '')

    with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
        with open(results_csvs_new+file+'.csv', 'w+', encoding="utf8") as csvoutput
        :
            writer = csv.writer(csvoutput, lineterminator='\n')
            reader = csv.reader(csvinput)

```

```

all = []
row = next(reader)

row.append('Bills')
all.append(row)

for row in reader:
    if re.search(r"(S\.\d{4})",row[12]) or re.search(r"(S\.\d{4})"
,row[12]) or re.search(r"(S\d{4})",row[12]) or re.search(r"(S \d{4})",row[12]) or re.se
arch(r"(H\.\R\.\.\d{4})",row[12]) or re.search(r"(HR \d{4})",row[12]) or re.search(r"(H
\.\R\.\d{4})",row[12]) or re.search(r"(HR\d{4})",row[12]):
        #if re.search(r"(.)*(S\.\d{4})*(S\.\d{4})*(S\d{4})*(S \d{4})*
(H\.\R\.\.\d{4})*(HR \d{4})*(H\.\R\.\d{4})*(HR\d{4})*(.)*$",row[12]):
            row.append('1')
            count += 1
            # print(count)
        else:
            row.append('0')

    all.append(row)
    #break
writer.writerows(all)

```

In []:

```

print ("No. of bills found : ")
print (count)

```

In [3]:

```
import requests
import os
import json
import xmltodict

import csv
import pandas as pd

months = ['01', '02', '03', '04', '05', '06', '07', '08', '09', '10', '11', '12']
years = ['1995', '1996', '1997', '1998', '1999', '2000', '2001', '2002', '2003', '2004',
        '2005', '2006', '2007', '2008', '2009', '2010', '2011', '2012']
committees = [102, 104, 106, 113, 115, 124, 128, 134, 138, 142, 156, 164, 173, 176, 182,
        184, 186, 192, 196, 242, 251, 305, 306, 308, 314, 316, 321, 330, 332, 336, 338, 344,
        358, 362, 380, 384, 388, 419, 432, 434, 435]
congresses = [104, 105, 106, 107, 108, 109, 110, 111, 112]

gpoShort = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_fr
om_server/gpo_tools/Extras/ChenJohnson_Agencies.csv"

df = pd.read_csv(gpoShort)
agencies = []

for i in (df['Agency']):
    agencies.append(i)

JK = []
UA = []
Parent = []

for i in (df['JK Code']):
    JK.append(i)
for i in (df['UA Code']):
    UA.append(i)
for i in (df['Parent UA Code']):
    Parent.append(i)
```


In [4]:

```
# CSV 1: Number of utterances made by the agency about a bill per month

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/CSV1.csv"

#for file in set(os.listdir(results_csvs)):
    #print (file)
    #print (set(os.listdir(results_csvs)))
    #file = file + '.csv'
    #if file in set(os.listdir(results_csvs)):
#file = file.replace('.csv','')

#with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    writer = csv.writer(csvoutput, lineterminator='\n')
    writer.writerow(["Date", "Committee", "Agency", "JK Code", "UA Code", "Parent UA Code"])

    for committee in committees:
        for month in months:
            for year in years:
                for i in range(len(agencies)):
                    row_temp = "=" + month + '-' + year + "\n", committee, agencies[i], JK[i], UA[i], Parent[i]
                    writer.writerow(row_temp)
```

In [2]:

```

# CSV 1: Number of utterances made by the agency about a bill per month

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

CSV1Dict = {}

utteranceCount = []

with open(CSV1, 'r', encoding="utf8") as csvinput2:
    CSV1reader = csv.reader(csvinput2)

    #all = []
    CSV1row = next(CSV1reader)

    #CSV1row.append('Number of utterances made by the agency about a bill per month')
    #all.append(CSV1row)

    for CSV1row in CSV1reader:
        count = 0

        CSV1RowDate = str(CSV1row[0])

```

```

CSV1RowDate = CSV1RowDate.replace('=', '')
CSV1RowDate = CSV1RowDate.replace('"', '')

CSV1key = CSV1RowDate+' '+ CSV1row[1] +' '+ CSV1row[2].strip()

CSV1Dict[CSV1key.strip()] = 0

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.keys())[0])

for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    with open(results_csvs+file+'.csv', 'r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            if row[27] == '1':

                date = row[13].split('-')[0]+'-'+row[13].split('-')

                indCSVkey = date +' '+ row[0] +' '+ row[21].strip()

                #print(indCSVkey)

                if indCSVkey.strip() in CSV1Dict.keys():
                    CSV1Dict[indCSVkey.strip()] += 1
                    #print(indCSVkey)

        #print(count)
        #utteranceCount.append(count)

```

 -
 StopIteration Traceback (most recent call last)

```

<ipython-input-2-1ae28d90f1c6> in <module>()
    47
    48     #all = []
--> 49     CSV1row = next(CSV1reader)
    50
    51     #CSV1row.append('Number of utterances made by the agency about
a bill per month')

```

StopIteration:

In [5]:

```

# To remove duplicate ent
CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
from more_itertools import unique_everseen
with open(CSV1, 'r') as f, open('2.csv', 'w') as out_file:
    out_file.writelines(unique_everseen(f))

```

In [3]:

```
print(len(CSV1Dict.keys()))  
print(list(CSV1Dict.values())[0:100000])  
print(CSV1Dict['06-1998 344 United States Postal Service'])
```

-
NameError Traceback (most recent call last)

<ipython-input-3-e85ae617f9f1> in <module>()

```
----> 1 print(len(CSV1Dict.keys()))  
      2 print(list(CSV1Dict.values())[0:100000])  
      3 print(CSV1Dict['06-1998 344 United States Postal Service'])
```

NameError: name 'CSV1Dict' is not defined

In [13]:

```

# CSV 1: Number of utterances made by the agency about a bill per month

CSV221 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from
_server/gpo_tools/CSV221.csv"

monthDict ={'01':'Jan',
            '02':'Feb',
            '03':'Mar',
            '04':'Apr',
            '05':'May',
            '06':'Jun',
            '07':'Jul',
            '08':'Aug',
            '09':'Sep',
            '10':'Oct',
            '11':'Nov',
            '12':'Dec'
            }

with open(CSV1,'r', encoding="utf8") as csvinput:
    with open(CSV221, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Number of utterances made by the agency about a bill pe
r month')
        all.append(row)

        for row in reader:
            CSV1RowDate = str(row[0])
            CSV1RowDate = CSV1RowDate.replace('=', '')
            CSV1RowDate = CSV1RowDate.replace('"', '')

            CSV1key = CSV1RowDate+' '+ row[1] + ' '+ row[2].strip()

            if CSV1key in CSV1Dict.keys():
                row.append(CSV1Dict[CSV1key])

            all.append(row)
            #break
        writer.writerows(all)

```

In [14]:

```
# Number of utterances made by the agency per month - CSV2

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

CSV1Dict = {}

utteranceCount = []

with open(CSV1, 'r', encoding="utf8") as csvinput2:
    CSV1reader = csv.reader(csvinput2)

    #all = []
    CSV1row = next(CSV1reader)

    #CSV1row.append('Number of utterances made by the agency about a bill per month')
    #all.append(CSV1row)

    for CSV1row in CSV1reader:
        count = 0

        CSV1RowDate = str(CSV1row[0])
```

```

CSV1RowDate = CSV1RowDate.replace('=', '')
CSV1RowDate = CSV1RowDate.replace('"', '')

CSV1key = CSV1RowDate+' '+ CSV1row[1] +' '+ CSV1row[2].strip()

CSV1Dict[CSV1key.strip()] = 0

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.keys())[0])

for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    with open(results_csvs+file+'.csv', 'r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            #if row[27] == '1':

                date = row[13].split('-')[0]+'-'+row[13].split('-')

[2]         indCSVkey = date +' '+ row[0] +' '+ row[21].strip

()

            #print(indCSVkey)

            if indCSVkey.strip() in CSV1Dict.keys():
                CSV1Dict[indCSVkey.strip()] += 1
                #print(indCSVkey)

        #print(count)
        #utteranceCount.append(count)

```

593352

01-1995 102 Broadcasting Board of Governors

In [15]:

```
# Number of utterances made by the agency per month - CSV2
```

```
CSV211 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from  
_server/gpo_tools/CSV211.csv"
```

```
monthDict ={'01':'Jan',  
            '02':'Feb',  
            '03':'Mar',  
            '04':'Apr',  
            '05':'May',  
            '06':'Jun',  
            '07':'Jul',  
            '08':'Aug',  
            '09':'Sep',  
            '10':'Oct',  
            '11':'Nov',  
            '12':'Dec'  
}
```

```
with open(CSV1,'r', encoding="utf8") as csvinput:  
    with open(CSV211, 'w+', encoding="utf8") as csvoutput:  
        writer = csv.writer(csvoutput, lineterminator='\n')  
        reader = csv.reader(csvinput)  
  
        all = []  
        row = next(reader)  
  
        row.append('Number of utterances made by the agency per month')  
        all.append(row)  
  
        for row in reader:  
            CSV1RowDate = str(row[0])  
            CSV1RowDate = CSV1RowDate.replace('=','')  
            CSV1RowDate = CSV1RowDate.replace('"','')  
  
            CSV1key = CSV1RowDate+' '+ row[1] + ' '+ row[2].strip()  
  
            if CSV1key in CSV1Dict.keys():  
                row.append(CSV1Dict[CSV1key])  
  
            all.append(row)  
            #break  
        writer.writerows(all)
```


In [17]:

```
# For each committee, need the number of total utterances per month - CSV3

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

CSV1Dict = {}

utteranceCount = []

with open(CSV1, 'r', encoding="utf8") as csvinput2:
    CSV1reader = csv.reader(csvinput2)

    #all = []
    CSV1row = next(CSV1reader)

    #CSV1row.append('Number of utterances made by the agency about a bill per month')
    #all.append(CSV1row)

    for CSV1row in CSV1reader:
        count = 0

        CSV1RowDate = str(CSV1row[0])
```

```

CSV1RowDate = CSV1RowDate.replace('=', '')
CSV1RowDate = CSV1RowDate.replace('"', '')

CSV1key = CSV1RowDate+' '+ CSV1row[1] #+' '+ CSV1row[2].strip()

CSV1Dict[CSV1key.strip()] = 0

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.keys())[0])

for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    with open(results_csvs+file+'.csv', 'r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            #if row[27] == '1':

                date = row[13].split('-')[0]+'-'+row[13].split('-')

                indCSVkey = date +' '+ row[0] #+' '+ row[21].strip

                #print(indCSVkey)

                if indCSVkey.strip() in CSV1Dict.keys():
                    CSV1Dict[indCSVkey.strip()] += 1
                    #print(indCSVkey)

        #print(count)
        #utteranceCount.append(count)

```

8856

01-1995 102

In [18]:

```
# For each committee, need the number of total utterances per month - CSV3

CSV2111 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_fro
m_server/gpo_tools/CSV2111.csv"

monthDict ={'01':'Jan',
            '02':'Feb',
            '03':'Mar',
            '04':'Apr',
            '05':'May',
            '06':'Jun',
            '07':'Jul',
            '08':'Aug',
            '09':'Sep',
            '10':'Oct',
            '11':'Nov',
            '12':'Dec'

            }

with open(CSV1,'r', encoding="utf8") as csvinput:
    with open(CSV2111, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Number of utterances made by the committees per month')
        all.append(row)

        for row in reader:
            CSV1RowDate = str(row[0])
            CSV1RowDate = CSV1RowDate.replace('=',' ')
            CSV1RowDate = CSV1RowDate.replace('"',' ')

            CSV1key = CSV1RowDate+' '+ row[1] #+' '+ row[2].strip()

            if CSV1key in CSV1Dict.keys():
                row.append(CSV1Dict[CSV1key])

            all.append(row)
            #break
        writer.writerows(all)
```

In [19]:

```
# For each agency, need the number of total utterances per month - CSV4

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

CSV1Dict = {}

utteranceCount = []

with open(CSV1, 'r', encoding="utf8") as csvinput2:
    CSV1reader = csv.reader(csvinput2)

    #all = []
    CSV1row = next(CSV1reader)

    #CSV1row.append('Number of utterances made by the agency about a bill per month')
    #all.append(CSV1row)

    for CSV1row in CSV1reader:
        count = 0

        CSV1RowDate = str(CSV1row[0])
```

```

CSV1RowDate = CSV1RowDate.replace('=', '')
CSV1RowDate = CSV1RowDate.replace('"', '')

CSV1key = CSV1RowDate+' '+ CSV1row[1].strip()

CSV1Dict[CSV1key.strip()] = 0

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.keys())[0])

for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    with open(results_csvs+file+'.csv', 'r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            #if row[27] == '1':

                date = row[13].split('-')[0]+'-'+row[13].split('-')

                indCSVkey = date + ' '+ row[21].strip()

                #print(indCSVkey)

                if indCSVkey.strip() in CSV1Dict.keys():
                    CSV1Dict[indCSVkey.strip()] += 1
                    #print(indCSVkey)

        #print(count)
        #utteranceCount.append(count)

```

14472

01-1995 Broadcasting Board of Governors

In [20]:

```
# For each agency, need the number of total utterances per month - CSV4
```

```
CSV4 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/CSV4.csv"
```

```
monthDict = {'01': 'Jan',
             '02': 'Feb',
             '03': 'Mar',
             '04': 'Apr',
             '05': 'May',
             '06': 'Jun',
             '07': 'Jul',
             '08': 'Aug',
             '09': 'Sep',
             '10': 'Oct',
             '11': 'Nov',
             '12': 'Dec'
            }

with open(CSV1, 'r', encoding="utf8") as csvinput:
    with open(CSV4, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Number of utterances made by the agencies per month')
        all.append(row)

        for row in reader:
            CSV1RowDate = str(row[0])
            CSV1RowDate = CSV1RowDate.replace('=', '')
            CSV1RowDate = CSV1RowDate.replace('"', '')

            CSV1key = CSV1RowDate + ' ' + row[1].strip()

            if CSV1key in CSV1Dict.keys():
                row.append(CSV1Dict[CSV1key])

            all.append(row)
            #break
        writer.writerows(all)
```

In [6]:

```
# Number of hearings made by the agency per month - CSV5

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

CSV1Dict = {}

utteranceCount = []

with open(CSV1, 'r', encoding="utf8") as csvinput2:
    CSV1reader = csv.reader(csvinput2)

    #all = []
    CSV1row = next(CSV1reader)

    #CSV1row.append('Number of utterances made by the agency about a bill per month')
    #all.append(CSV1row)

    for CSV1row in CSV1reader:
        count = 0

        CSV1RowDate = str(CSV1row[0])
```

```

CSV1RowDate = CSV1RowDate.replace('=', '')
CSV1RowDate = CSV1RowDate.replace('"', '')

CSV1key = CSV1RowDate+' '+ CSV1row[1] +' '+ CSV1row[2].strip()

CSV1Dict[CSV1key.strip()] = 0

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.keys())[0])

hearingsSet = set()
for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    hearingsSet.clear()
    with open(results_csvs+file+'.csv', 'r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            #if row[27] == '1':

                date = row[13].split('-')[0]+'-'+row[13].split('-')

[2]         indCSVkey = date +' '+ row[0] +' '+ row[21].strip

()

                hearingsSet.add(indCSVkey)
                #print(indCSVkey)

    for i in hearingsSet:
        if i.strip() in CSV1Dict.keys():
            CSV1Dict[i.strip()] += 1
            #print(indCSVkey)

    #print(count)
    #utteranceCount.append(count)

print(len(CSV1Dict.keys()))
print(list(CSV1Dict.values())[0])

```

```

593352
01-1995 102 Broadcasting Board of Governors
593352
0

```


In [7]:

```

# Number of hearings made by the agency per month - CSV5

CSV5 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV5.csv"

monthDict ={'01':'Jan',
            '02':'Feb',
            '03':'Mar',
            '04':'Apr',
            '05':'May',
            '06':'Jun',
            '07':'Jul',
            '08':'Aug',
            '09':'Sep',
            '10':'Oct',
            '11':'Nov',
            '12':'Dec'

            }

with open(CSV1,'r', encoding="utf8") as csvinput:
    with open(CSV5, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Number of hearings made by the agency per month')
        all.append(row)

        for row in reader:
            CSV1RowDate = str(row[0])
            CSV1RowDate = CSV1RowDate.replace('=', '')
            CSV1RowDate = CSV1RowDate.replace("'", '')

            CSV1key = CSV1RowDate+' '+ row[1] + ' '+ row[2].strip()

            if CSV1key in CSV1Dict.keys():
                row.append(CSV1Dict[CSV1key])

            all.append(row)
            #break
        writer.writerows(all)

```

In [24]:

```

# Finding gender based on names

import re
import sys
import csv

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

CSV1 = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_s
erver/gpo_tools/CSV1.csv"
results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearing
s_from_server/gpo_tools/results_csvs/"

namesDict = {}

for file in set(os.listdir(results_csvs)):

    file = file.replace('.csv', '')

    hearingsSet.clear()
    with open(results_csvs+file+'.csv','r', encoding="utf8") as csvinput:
        # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
        #     writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)
        row = next(reader)
        for row in reader:
            if row[17].strip() != 'NA' or row[17].strip() != '-':
                namesDict[row[17]] = 'M/F'
                #print(indCSVkey)

```

```
print(len(namesDict.keys()))  
print(list(namesDict.values())[0])
```

56835

M/F

In [38]:

```
# Finding gender based on names
```

```
import gender_guesser.detector as gender
```

```
d = gender.Detector()
```

```
for i in range(30):
```

```
    print((list(namesDict.keys())[i]) + " : "+ d.get_gender(list(namesDict.keys())[i]))
```

```
    #print('\n')
```

```
print(d.get_gender(u"Mainzer"))
```

```
print(d.get_gender(u"Bob"))
```

```
hastings, richard doc : unknown
```

```
defazio, peter a : unknown
```

```
Mainzer Elliot : unknown
```

```
Kem John : unknown
```

```
Eichenberger Kathy : unknown
```

```
Brigham Kathryn : unknown
```

```
Corwin Scott : unknown
```

```
Webb Tony : unknown
```

```
Crinklaw Rick : unknown
```

```
Reimann Ron : unknown
```

```
Haller Greg : unknown
```

```
Amos Paul : unknown
```

```
McCart Wes : unknown
```

```
Spencer Bachus : unknown
```

```
Carolyn B. Maloney : unknown
```

```
Jeb Hensarling : unknown
```

```
David Scott : unknown
```

```
royce, ed : unknown
```

```
Barney Frank : unknown
```

```
Sean P. Duffy : unknown
```

```
Stephen F. Lynch : unknown
```

```
canseco, francisco : unknown
```

```
Brad Miller : unknown
```

```
Walter B. Jones Jr. : unknown
```

```
slaughter, louise m : unknown
```

```
walz, timothy j : unknown
```

```
Judy Biggert : unknown
```

```
Donald A. Manzullo : unknown
```

```
Maxine Waters : unknown
```

```
Bill Posey : unknown
```

```
unknown
```

```
male
```

In [27]:

```
import nltk  
nltk.download('names')
```

```
[nltk_data] Downloading package names to  
[nltk_data] C:\Users\RAHUL\AppData\Roaming\nltk_data...  
[nltk_data] Unzipping corpora\names.zip.
```

Out[27]:

True

In [39]:

```
# Finding gender based on names

import random
from nltk.corpus import names
import nltk

def gender_features(word):
    return {'last_letter':word[-1]}

# preparing a list of examples and corresponding class labels.
labeled_names = ([ (name, 'male') for name in names.words('male.txt') ] +
                  [ (name, 'female') for name in names.words('female.txt') ])

random.shuffle(labeled_names)

# we use the feature extractor to process the names data.
featuresets = [(gender_features(n), gender)
                for (n, gender) in labeled_names]

# Divide the resulting list of feature
# sets into a training set and a test set.
train_set, test_set = featuresets[5:], featuresets[:5]

# The training set is used to
# train a new "naive Bayes" classifier.
classifier = nltk.NaiveBayesClassifier.train(train_set)

for i in range(30):
    print((list(namesDict.keys())[i]) + " : " + classifier.classify(gender_features(list(
namesDict.keys())[i])))
    #print('\n')

print(classifier.classify(gender_features('Bob')))
```

hastings, richard doc : male
defazio, peter a : female
Mainzer Elliot : male
Kem John : male
Eichenberger Kathy : female
Brigham Kathryn : male
Corwin Scott : male
Webb Tony : female
Crinklaw Rick : male
Reimann Ron : male
Haller Greg : male
Amos Paul : male
McCart Wes : male
Spencer Bachus : male
Carolyn B. Maloney : female
Jeb Hensarling : male
David Scott : male
royce, ed : male
Barney Frank : male
Sean P. Duffy : female
Stephen F. Lynch : female
canseco, francisco : male
Brad Miller : male
Walter B. Jones Jr. : female
slaughter, louise m : male
walz, timothy j : male
Judy Biggert : male
Donald A. Manzullo : male
Maxine Waters : male
Bill Posey : female
male

In [40]:

```
print(len(namesDict.keys()))  
print(list(namesDict.values())[0])
```

56835

M/F

In [41]:

```
fout = "namesDict.txt"  
fo = open(fout, "w")  
  
for k, v in namesDict.items():  
    fo.write(str(k) + '\n')  
  
fo.close()
```

In [43]:

```
count = 0

for k, v in namesDict.items():
    if (str(k).find(',') != -1 ):
        count += 1

print (count)
```

1122

In [1]:

```
# Metadata
# subCommittee extraction

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```



```

df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('subCommittee')
        all.append(row)

        for row in reader:
            try:

                file = row[6] + ".json"

                with open(APIs+file) as data_file:
                    jsonObj = json.load(data_file)

                    if (jsonObj["mods"]["extension"][2]["congressCommittee"]["subCommittee"]["name"]):
                        subCommittee = jsonObj["mods"]["extension"][2]["congressCommittee"]["subCommittee"]["name"]
                        row.append(subCommittee)
                        #print(subCommittee)
                    else:
                        row.append('-')

            except:
                row.append("-")

            all.append(row)

        #except:
        #    writer.writerow(all)
        #    continue
        writer.writerow(all)

```

C:\Users\RAHUL\Anaconda3\lib\site-packages\fuzzywuzzy\fuzz.py:11: UserWarning: Using slow pure-python SequenceMatcher. Install python-Levenshtein to remove this warning

warnings.warn('Using slow pure-python SequenceMatcher. Install python-Levenshtein to remove this warning')

In [4]:

```
# Metadata
# Column: "Committee member count"

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```
df1 = pd.read_csv(metadata_results)
#print(df1['filename'])

sample_jackets = ['CHRG-115hhrg27211']
count = 0

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Committee member count')
        all.append(row)

        for row in reader:
            count = len(row[14].split('\n'))

            row.append(count)

            all.append(row)

        #except:
        #    writer.writerow(all)
        #    continue
        writer.writerow(all)
```

In [31]:

```
# Metadata
# Column: "Denominator count"

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

House = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
server/gpo_tools/Extras/house_assignments_103-115-3.csv"
Senate = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
_server/gpo_tools/Extras/senate_assignments_103-115-3.csv"

CongCom = {}

with open(House, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        if (row[0]+'_'+row[1] in CongCom.keys()):
            CongCom[row[0]+'_'+row[1]] += 1
        else:
            CongCom[row[0]+'_'+row[1]] = 1

with open(Senate, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        if (row[0]+'_'+row[1] in CongCom.keys()):
            CongCom[row[0]+'_'+row[1]] += 1
        else:
            CongCom[row[0]+'_'+row[1]] = 1

#print(CongCom)

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Denominator count')
        all.append(row)

        for row in reader:
            if (row[2].replace("th", "")+'_'+row[3]) in CongCom.keys():
                count = CongCom[(row[2].replace("th", "")+'_'+row[3])]
            else:
                count = '-'

            row.append(count)
            all.append(row)

        #except:
        #     writer.writerow(all)
        #     continue
        writer.writerow(all)

```

In [32]:

```
# Metadata
# Column: "Party count"

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

House = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
server/gpo_tools/Extras/house_assignments_103-115-3.csv"
Senate = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
_server/gpo_tools/Extras/senate_assignments_103-115-3.csv"

CongCom = {}

with open(House, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        if (row[0]+'::'+row[1]+'::'+row[6] in CongCom.keys()):
            CongCom[row[0]+'::'+row[1]+'::'+row[6]] += 1
        else:
            CongCom[row[0]+'::'+row[1]+'::'+row[6]] = 1

with open(Senate, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        if (row[0]+'::'+row[1]+'::'+row[6] in CongCom.keys()):
            CongCom[row[0]+'::'+row[1]+'::'+row[6]] += 1
        else:
            CongCom[row[0]+'::'+row[1]+'::'+row[6]] = 1

#print(CongCom)

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Party count(100:200:328:999:9999)')
        all.append(row)

        for row in reader:
            if (row[2].replace("th", "")+'::'+row[3]+'::100') in CongCom.k
            count100 = CongCom[(row[2].replace("th", "")+'::'+row[3]+'
            ':100')]
            else:
                count100 = '-'

            if (row[2].replace("th", "")+'::'+row[3]+'::200') in CongCom.k
            count200 = CongCom[(row[2].replace("th", "")+'::'+row[3]+'
            ':200')]
            else:
                count200 = '-'

            if (row[2].replace("th", "")+'::'+row[3]+'::328') in CongCom.k
            count328 = CongCom[(row[2].replace("th", "")+'::'+row[3]+'
            ':328')]
            else:
                count328 = '-'

            row.append(count100)
            row.append(count200)
            row.append(count328)
            all.append(row)

        writer.writerow(all)

```

```

count328 = CongCom[(row[2].replace("th","")+':'+row[3]+
':328')]

else:
    count328 = '-'

if (row[2].replace("th","")+':'+row[3]+':999') in CongCom.k
keys():
    count999 = CongCom[(row[2].replace("th","")+':'+row[3]+
':999')]

else:
    count999 = '-'

if (row[2].replace("th","")+':'+row[3]+':9999') in CongCom.
keys():
    count9999 = CongCom[(row[2].replace("th","")+':'+row[3]
+':9999')]

else:
    count9999 = '-'

    temp = "=" + str(count100) + ":" + str(count200) + ":" +
str(count328) + ":" + str(count999) + ":" + str(count9999) + "\n"

    row.append( temp)
    all.append(row)

except:
    # writer.writerow(all)
    # continue
writer.writerow(all)

```

In [26]:

```

#Finding unique party codes
PartyCodes = {}

with open(House,'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        if(row[6] in PartyCodes.keys()):
            PartyCodes[row[6]] += 1
        else:
            PartyCodes[row[6]] = 0

print(PartyCodes.keys())

dict_keys(['200', '100', '328', '999', '9999'])

```


In [39]:

```
# Metadata
# Column: "Party & Committee info:"

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

House = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
server/gpo_tools/Extras/house_assignments_103-115-3.csv"
Senate = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
server/gpo_tools/Extras/senate_assignments_103-115-3.csv"

PartyCom = {}

with open(House, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        PartyCom[row[0]+row[1]+row[3].lower().strip()] = row[6]
+':'+row[9]+':'+row[10]

with open(Senate, 'r', encoding="utf8") as csvinput:
    # with open(CSV1, 'w+', encoding="utf8") as csvoutput:
    #     writer = csv.writer(csvoutput, lineterminator='\n')
    reader = csv.reader(csvinput)
    row = next(reader)
    for row in reader:
        PartyCom[row[0]+row[1]+row[3].lower().strip()] = row[6]
+':'+row[10]+':'+row[11]

#print(CongCom)

with open(metadata_results, 'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Party & Committee info(Party:Senior Party Member:Committee Seniority)')
        all.append(row)

        for row in reader:
            temp = []
            for name in row[14].split('\n'):
                if row[2].replace("th", "")+row[3]+name.split(' : ')[0].lower().strip().replace(';','') in PartyCom.keys():
                    temp.append( PartyCom[row[2].replace("th", "")+row[3]+name.split(' : ')[0].lower().strip().replace(';','')] )
                else:
                    temp.append( '-'+'-'+'-' )

            row.append("\n".join(temp))
            all.append(row)

        #except:
        #     writer.writerow(all)
        #     continue
        writer.writerow(all)

```

In [45]:

```
# Metadata
# Column: "Expertise"

import json
from pprint import pprint

import sys
import csv

from fuzzywuzzy import fuzz
from fuzzywuzzy import process

#csv.field_size_limit(sys.maxsize)

maxInt = sys.maxsize
decrement = True

while decrement:
    # decrease the maxInt value by factor 10
    # as long as the OverflowError occurs.

    decrement = False
    try:
        csv.field_size_limit(maxInt)
    except OverflowError:
        maxInt = int(maxInt/10)
        decrement = True

import os

import math

import requests
import xml.etree.ElementTree as ET

import json
import xmltodict

import csv
import pandas as pd

metadata_results = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results.csv"
metadata_results_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/metadata_results_new.csv"

results_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs/"
results_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/results_csvs_new/"

sample_csvs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs/"
sample_csvs_new = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/sample_csvs_new/"

APIs = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_server/gpo_tools/APIs/"
```

```

House = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
server/gpo_tools/Extras/house_assignments_103-115-3.csv"
Senate = "D:/USC/RA NLP/Hearing data/congressional_hearings/congressional_hearings_from_
_server/gpo_tools/Extras/senate_assignments_103-115-3.csv"

expertise={

    'A.A.' : 'Associate of Arts',
        'A.S.' : 'Associate of Science',
        'A.A.S.' : 'Associate of Applied Science',
        'ADN' : 'Associates Degree in Nursing',
        'B.A.' : 'Bachelor of Arts',
        'B.S.' : 'Bachelor of Science',
        'B.E.' : 'Bachelor of Engineering',
        'M.A.' : 'Master of Arts',
        'M.S.' : 'Master of Science',
        'MBA' : 'Master of Business Administration',
        'M.Ed.' : 'Master of Education',
        'Ph.D.' : 'Doctor of Philosophy',
        'DNP' : 'Doctor of Nursing Practice',
        'Ed.D.' : 'Doctor of Education',
    'J.D.' : 'Juris Doctorate, a law degree',
        'M.D.' : 'Medical Doctor, a physicians degree',
        'D.D.S.' : 'Doctor of Dental Surgery, a dentistry degree',
    'Pharm.D.' : 'Doctor of Pharmacy , a pharmaceutical medicine degree'

}

with open(metadata_results,'r', encoding="utf8") as csvinput:
    with open(metadata_results_new, 'w+', encoding="utf8") as csvoutput:
        writer = csv.writer(csvoutput, lineterminator='\n')
        reader = csv.reader(csvinput)

        all = []
        row = next(reader)

        row.append('Expertise')
        all.append(row)

        for row in reader:
            temp = []
            for name in row[13].split('\n'):
                done = 0
                for i in name.split():
                    #print (i)
                    if i.strip() in expertise.keys():
                        temp.append(i+' : '+expertise[i])
                        done = 1
                    break
                # print(i)

            if done == 0:
                temp.append('-')

            row.append("\n".join(temp))
            all.append(row)

        #break

    #except:

```

```
#    writer.writerows(all)
#    continue
writer.writerows(all)
```

In []: