# Aim:

In this analysis we are going to perform entity extraction on the Reuters corpus and construct entity profiles for persons, organizations and locations.

**Time required to run code: The entire Notebook finishes execution in about *600 seconds or 10 mins*, this is due to creating Spacy Doc objects for each document for the Reuters Dataset.**

# Index

---

# 1. Libraries needed to run notebook

# 2. Approach to solving the problem

# 3. Code

# 4. Result Analysis

## a. Top entities for all Reuters document

## b. Most Frequent Terms for popular entities

# 5. Extra Credit

## a. Approach to solving the extra credit problem

# 6. Future Scope

---

# 1. Libraries needed to run notebook:

The following libraries are required to run this Notebook:

1. NLTK
2. Spacy
3. Matplotlib
4. wordcloud

# 2. Approach for solving the Problem:

## a. Preprocessing:

The preprocessing code can be found in Step 2 of the code section.

Since each Reuters document contains new line character '\n', we remove them, since it causes hindrance for Spacy's entity extraction. To remove new lines we have used **re.sub(r"\n", "", doc_text)**. We have also replaced any spaces which occur more than once with a single space. The following code snippet was used to remove more than 1 spaces: **re.sub(r"\s+", " ", corpus_clean)**. Thus finally we get a paragraph from the Reuters document.

In the Reuters document there was an incompatible unicode character "&lt", which was replaced with "<".

Since Reuters documents contain information about Company's stock market ticker, eg: Google has < GOOGL >, we removed these tickers by using **re.sub(r"<.*>", " ", corpus_clean)**.

We also remove all special characters except for period(.), comma(,) and apostrophe('), for this we use **re.sub(r"[^A-Za-z0-9.,']", " ", corpus_clean)**.

## b. Extracting entities from a single Reuters Document:

The **extract_entities(doc_id, doc_text)** function in Step 3 of the code section performs the entity extraction after the preprocessing is complete.

We keep track of different entities like person, location and organization in different dictionaries. Using Spacy's .ents attribute we get the entity type and filter it based on "PERSON", "GPE" and "ORG". Each respective dictionary has the following structure, the **key** is the entity text, which is in the lowercase format for avoiding different forms of entity name and the **value** is a dictionary in which the **key** is document id and **value** is a list of sentences where the entity occurs.

eg: The person dictionary will be as follows:

```
{"Sheldon":
    {"doc_id_1":
        ["My name is Sheldon" ,
         "Sheldon likes the color red.",
         .....
        ]
    },
 ......
}
```

The same logic is repeated for location and organization.

## c. Storing document entities in a combined dictionary:

The Step 4 in the code section is used to store all document entities in a combined dictionary. After extracting entities from a single Reuters document we add that dictionary into a combined dictionary for which the **key** is the entity name and the **value** is a list in which the *0 index* contains a dictionary in which the **key** is document id and **value** is list of all the sentences in which the entity occurs and the *index 1* contains the *count* of the sentences in which the entity occurs across all the documents.

eg: The combined person dictionary will be as follows:

```
{"Sheldon": [
        {
            "doc_id_1":
                ["My name is Sheldon" ,
                 "Sheldon likes the color red."
                ],
            "doc_id_2":
                ["My name is also Sheldon" ,
                 "Sheldon likes the color blue."
                ]
        },
        4
    ]
}
```

The same logic is repeated for location and organization.

## d. Finding most popular entities:

The Step 5 of code section is used to get the most popular entities in the entire Reuters dataset. We have sorted each combined dictionary in descending order based on the number of times they occur in the Reuters dataset and then we return the top 500 entities.

## e. Finding the words which most frequently occur with most popular entities:

The **find_frequent_words_with_entity(top_entities, combined_entities)** function from Step 7 of the code section is used to get the words which most frequently occur with the most popular entities. We find the popular **nouns, verbs and adjectives** for a given entity.

We combine all the sentences for a popular entity into a single paragraph and also remove the entity name itself to ignore it as a frequent term. We then create a Spacy Doc object using this paragraph and filter out the stop words, numbers, dates to create a **lemmatized list** for noun, verb and adjectives individually.

We then store the top 10 most frequent tokens for noun, verb and adjective; and return a dictionary of the following format:

```
{"entity_name":[
            [frequent nouns],
            [frequent verbs],
            [frequent adjectives]
        ]
}
```

This logic is used for person, location and organization.

---

# 3. Code:

## Step 1) Import necessary libraries:

```
In [1]:    1  from collections import Counter
           2  from nltk.corpus import reuters
           3  from wordcloud import WordCloud
           4  import matplotlib.pyplot as plt
           5  import spacy
           6  import re
```

We are working with the Large Wikipedia English trained model, we can install it using:

```
python -m spacy download en_core_web_lg
```

```
In [2]:    1  nlp = spacy.load("en_core_web_lg")
```

## Step 2) Preprocessing for Reuters Document:

```
In [3]:    1  def preprocess_doc_text(doc_text):
           2      # remove all new line characters from corpus
           3      corpus_clean = re.sub(r"\n", "", doc_text)
           4
           5      # replace all more than 1 spaces with single space
           6      corpus_clean = re.sub(r"\s+", " ", corpus_clean)
           7
           8      # fix weird unicode in text
           9      corpus_clean = re.sub(r"&lt;", "<", corpus_clean)
          10
          11      # remove the market ticker name of the company
          12      corpus_clean = re.sub(r"<.*>", " ", corpus_clean)
          13
          14      # remove all special characters from the corpus, except for period, commaand apost
          15      corpus_clean = re.sub(r"[^A-Za-z0-9.,']", " ", corpus_clean)
          16
          17      return corpus_clean
```

## Step 3) Function to extract the entity, document id, and relevant sentence text from the input:

```
In [4]:   1  def extract_entities(doc_id, doc_text):
          2
          3      # create spacy doc object
          4      analyzed_doc = nlp(preprocess_doc_text(doc_text))
          5
          6      doc_persons = {}
          7      doc_organizations = {}
          8      doc_locations = {}
          9
         10      # analyzing entities for a document
         11      for entity in analyzed_doc.ents:
         12          entity_text = entity.text.lower().strip()
         13          sentence = entity.sent.text.strip()
         14
         15          if entity_text != "":
         16
         17              # for persons
         18              if entity.label_ == "PERSON":
         19                  if entity_text in doc_persons.keys():
         20                      # present
         21                      doc_persons[entity_text][doc_id].append(sentence)
         22                  else:
         23                      # not present
         24                      doc_persons[entity_text] = {doc_id: [sentence]}
         25
         26              # for locations
         27              if entity.label_ == "GPE":
         28                  if entity_text in doc_locations.keys():
         29                      # present
         30                      doc_locations[entity_text][doc_id].append(sentence)
         31                  else:
         32                      # not present
         33                      doc_locations[entity_text] = {doc_id: [sentence]}
         34
         35              # for organizations
         36              if entity.label_ == "ORG":
         37                  if entity_text in doc_organizations.keys():
         38                      # present
         39                      doc_organizations[entity_text][doc_id].append(sentence)
         40                  else:
         41                      # not present
         42                      doc_organizations[entity_text] = {doc_id: [sentence]}
         43
         44      return doc_persons, doc_organizations, doc_locations
```

## Step 4) Calling extract_entities on all Reuters documents:

```python
num_docs = len(reuters.fileids())

combined_persons = {}
combined_organizations = {}
combined_locations = {}

for doc_id in reuters.fileids()[:num_docs]:

    persons, organizations, locations = extract_entities(
        doc_id, reuters.open(doc_id).read()
    )

    for person in persons.keys():
        if person in combined_persons.keys():
            # present
            combined_persons[person][0][doc_id] = persons[person][doc_id]
            # update frequency of entity occuring across all documents
            combined_persons[person][1] += len(persons[person][doc_id])
        else:
            # not present
            combined_persons[person] = [persons[person], len(persons[person][doc_id])]

    for organization in organizations.keys():
        if organization in combined_organizations.keys():
            # present
            combined_organizations[organization][0][doc_id] = organizations[organizati
            combined_organizations[organization][1] += len(organizations[organization]
        else:
            # not present
            combined_organizations[organization] = [
                organizations[organization],
                len(organizations[organization][doc_id]),
            ]

    for location in locations.keys():
        if location in combined_locations.keys():
            # present
            combined_locations[location][0][doc_id] = locations[location][doc_id]
            combined_locations[location][1] += len(locations[location][doc_id])
        else:
            # not present
            combined_locations[location] = [
                locations[location],
                len(locations[location][doc_id]),
            ]
```

## Step 5) To get the most popular entities based on number of times they occur in entire Reuters dataset:

```
In [6]:  1  def find_most_popular_entities(entity_dictionary):
         2      list_of_dictionary_keys_with_most_mentions = []
         3
         4      for k, v in sorted(entity_dictionary.items(), key=lambda item: -item[1][1]):
         5          # store the entity name and its frequency
         6          list_of_dictionary_keys_with_most_mentions.append((k, v[1]))
         7
         8      # return the top 500 entities
         9      return list_of_dictionary_keys_with_most_mentions[:500]
```

## Step 6) Invoke top entity mention finder:

```
In [7]:  1  top_persons = find_most_popular_entities(combined_persons)
         2  top_organizations = find_most_popular_entities(combined_organizations)
         3  top_locations = find_most_popular_entities(combined_locations)
```

## Step 7) Analyzing most popular entities to determine what words they most frequently occur with:

```python
In [8]:  1  def find_frequent_words_with_entity(top_entities, combined_entities):
         2      most_popular_terms = {}
         3
         4      # finally, now find the most frequent tokens associated with the entities
         5      for entity, frequency in top_entities:
         6          # using the top_persons list find get the values in combined_persons
         7          entity_details, entity_count = combined_entities[entity]
         8
         9          # store all the sentences across documents into a common list
        10          common_sent = []
        11          for statements in list(entity_details.values()):
        12              common_sent.extend(statements)
        13
        14          # get all the sentences into a paragraph
        15          joined_sentence = " ".join(common_sent)
        16
        17          # we remove the entity name itself, to avoid counting it
        18          sentence_processed = re.sub(entity, "", joined_sentence, flags=re.IGNORECASE)
        19
        20          # lower case the sentence for frequency counting
        21          final_paragraph = nlp(sentence_processed.lower())
        22
        23          # store the verb, noun and adjective
        24          noun_words = []
        25          verb_words = []
        26          adjective_words = []
        27
        28          for token in final_paragraph:
        29              # remove all stop words and symbols
        30              # remove all numbers and dates
        31              # lemmatize all words
        32              if (
        33                  token.is_stop != True
        34                  and token.is_punct != True
        35                  and token.like_num != True
        36                  and token.text.strip() != ""
        37              ):
        38                  if token.pos_ == "NOUN" or token.pos_ == "PROPN":
        39                      noun_words.append(token.lemma_)
        40                  elif token.pos_ == "ADJ":
        41                      adjective_words.append(token.lemma_)
        42                  elif token.pos_ == "VERB":
        43                      verb_words.append(token.lemma_)
        44
        45          # find frequency
        46          noun_word_freq = Counter(noun_words)
        47          verb_word_freq = Counter(verb_words)
        48          adjective_word_freq = Counter(adjective_words)
        49
        50          # return the top 10 tokens
        51          noun_common_words = noun_word_freq.most_common(10)
        52          verb_common_words = verb_word_freq.most_common(10)
        53          adjective_common_words = adjective_word_freq.most_common(10)
        54
        55          # fill this dictionary with all the words in the context of the entity
        56          most_popular_terms[entity] = [
        57              # 0 = nouns and proper nouns
        58              [frequent_words[0] for frequent_words in noun_common_words],
        59              # 1 verbs
        60              [frequent_words[0] for frequent_words in verb_common_words],
        61              # 2 adjectives
```

```
62                    [frequent_words[0] for frequent_words in adjective_common_words],
63            ]
64
65        return most_popular_terms
```

In [9]:
```
1  # finally, now find the most frequent tokens associated with the entities
2  person_most_popular_terms = find_frequent_words_with_entity(
3      top_persons, combined_persons
4  )
5
6  organization_most_popular_terms = find_frequent_words_with_entity(
7      top_organizations, combined_organizations
8  )
9
10 location_most_popular_terms = find_frequent_words_with_entity(
11     top_locations, combined_locations
12 )
```

## Step 8) Function which creates word cloud based on frequency of terms

In [10]:
```
1  def create_word_cloud(entity_with_frequencies):
2      wordcloud = WordCloud(
3          background_color="white", width=1000, height=500
4      ).generate_from_frequencies(dict(entity_with_frequencies))
5
6      plt.figure(figsize=(13, 13), facecolor=None)
7      plt.imshow(wordcloud)
8      plt.axis("off")
9      plt.tight_layout(pad=0)
10
11     plt.show()
```
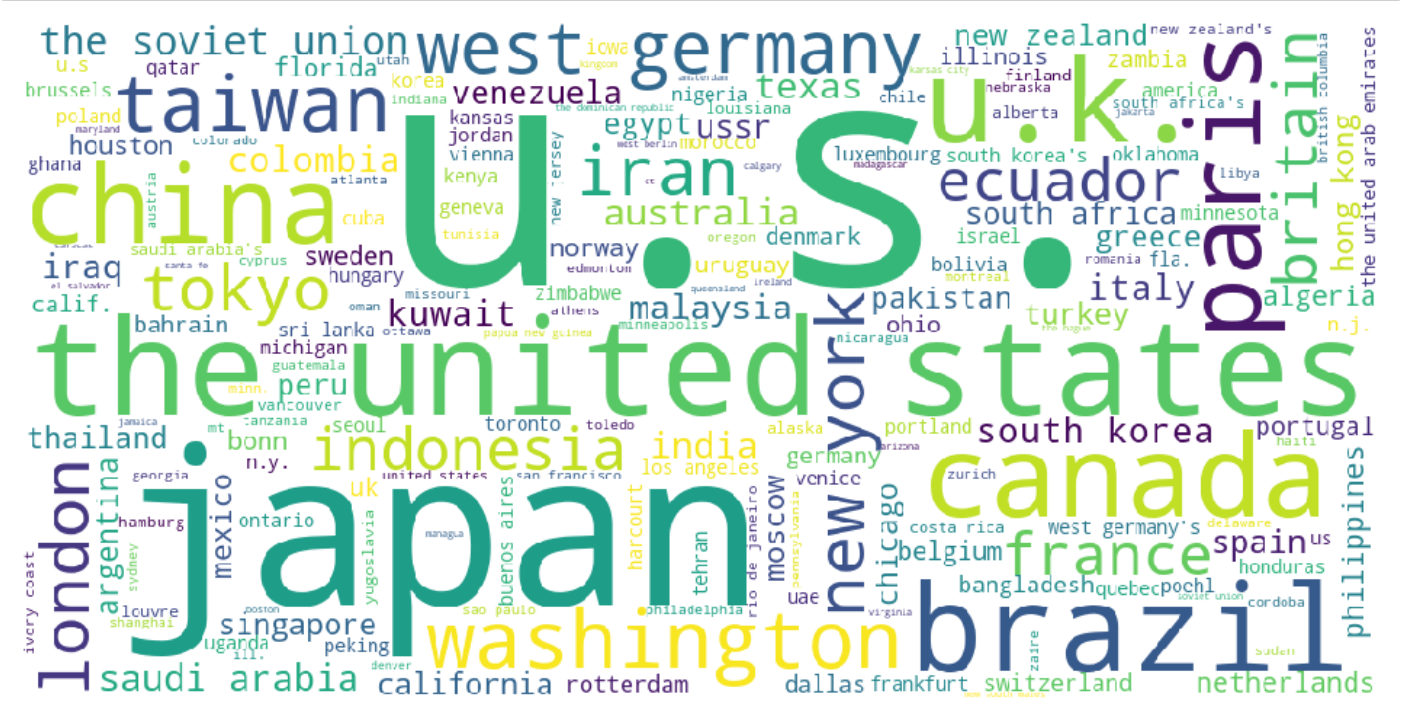
## 4. Result Analysis:

### a. Top entities for all Reuters document

We have found the top 500 entities across all documents for person, location and organization entity labels. The following word cloud represents the results:

i. Top 500 frequently occuring Persons:

ii. Top 500 frequently occuring Locations:

```
In [12]:   1  create_word_cloud(top_locations)
```



iii. Top 500 frequently occuring Organizations:

```
In [13]:   1  create_word_cloud(top_organizations)
```



## b. Most Frequent Terms for popular entities

We discuss the most frequent terms for the top 3 persons, locations and organizations.

The top 3 Persons are:

```
In [14]:    1  top_persons[:3]
```

Out[14]: [('reagan', 382), ('baker', 275), ('vs', 165)]

And the most frequently occurring terms with these entities are given below.

**Note:** Inside the nested list, the first list is for NOUNS, the second list is for VERBS and third list is for ADJECTIVES

```
In [15]:  1  for entity_name, count in top_persons[:3]:
          2      print("------------------------------------------------")
          3      print("The entity name is:", entity_name)
          4      print()
          5      print(
          6          "Most popular Nouns associate are: ", person_most_popular_terms[entity_name][0
          7      )
          8      print()
          9      print(
         10          "Most popular Verbs associate are: ", person_most_popular_terms[entity_name][1
         11      )
         12      print()
         13      print(
         14          "Most popular Adjectives associate are: ",
         15          person_most_popular_terms[entity_name][2],
         16      )
         17      print()
```

```
------------------------------------------------
The entity name is: reagan

Most popular Nouns associate are:  ['u.s', 'president', 'trade', 'administration', 'japa
n', 'bill', 'house', 'oil', 'official', 'agreement']

Most popular Verbs associate are:  ['say', 'impose', 'urge', 'retaliate', 'require', 'ope
n', 'announce', 'propose', 'protect', 'include']

Most popular Adjectives associate are:  ['japanese', 'foreign', 'economic', 'unfair', 'ne
w', 'american', 'protectionist', 'free', 'strong', 'domestic']

------------------------------------------------
The entity name is: baker

Most popular Nouns associate are:  ['treasury', 'u.s', 'rate', 'hughes', 'currency', 'dol
lar', 'merger', 'agreement', 'meeting', 'policy']

Most popular Verbs associate are:  ['say', 'see', 'tell', 'agree', 'decline', 'go', 'thin
k', 'comment', 'propose', 'meet']

Most popular Adjectives associate are:  ['german', 'west', 'economic', 'international',
'monetary', 'japanese', 'current', 'consistent', 'great', 'short']

------------------------------------------------
The entity name is: vs

Most popular Nouns associate are:  ['ct', 'shr', 'net', 'cts', 'mln', 'qtr', 'sale', 'cor
p', 'year', 'div']

Most popular Verbs associate are:  ['end', 'dilute', 'set', 'vote', 'give', 'make', 'adju
st', 'pre', 'raise']

Most popular Adjectives associate are:  ['net', '4th', '2nd', 'prior', '3rd', '1st', 'con
tinental', 'primary', 'quarterly', 'payable']


The top 3 Locations are:
```

```
In [16]:    1  top_locations[:3]
```

Out[16]: [('u.s.', 4289), ('japan', 1530), ('the united states', 580)]

And the most frequently occurring terms with these entities are given below.

**Note:** Inside the nested list, the first list is for NOUNS, the second list is for VERBS and third list is for ADJECTIVES

```
In [17]:  1  for entity_name, count in top_locations[:3]:
          2      print("------------------------------------------------")
          3      print("The entity name is:", entity_name)
          4      print()
          5      print(
          6          "Most popular Nouns associate are: ",
          7          location_most_popular_terms[entity_name][0],
          8      )
          9      print()
         10      print(
         11          "Most popular Verbs associate are: ",
         12          location_most_popular_terms[entity_name][1],
         13      )
         14      print()
         15      print(
         16          "Most popular Adjectives associate are: ",
         17          location_most_popular_terms[entity_name][2],
         18      )
         19      print()
```

```
------------------------------------------------
The entity name is: u.s.

Most popular Nouns associate are:  ['trade', 'japan', 'market', 'pct', 'year', 'dlr', 'do
llar', 'export', 'oil', 'official']

Most popular Verbs associate are:  ['say', 'rise', 'expect', 'tell', 'fall', 'buy', 'incr
ease', 'report', 'add', 'sell']

Most popular Adjectives associate are:  ['japanese', 'foreign', 'high', 'major', 'canadia
n', 'economic', 'new', 'iranian', 'low', 'large']

------------------------------------------------
The entity name is: japan

Most popular Nouns associate are:  ['u.s', 'trade', 'market', 'pct', 'year', 'official',
'surplus', 'united', 'states', 'import']

Most popular Verbs associate are:  ['say', 'cut', 'open', 'tell', 'rise', 'buy', 'reduc
e', 'increase', 'fall', 'continue']

Most popular Adjectives associate are:  ['ese', 'economic', 'domestic', 'foreign', 'larg
e', 'major', 'current', 'high', 'new', 'international']

------------------------------------------------
The entity name is: the united states

Most popular Nouns associate are:  ['trade', 'japan', 'u.s', 'official', 'country', 'oi
l', 'market', 'agreement', 'surplus', 'import']

Most popular Verbs associate are:  ['say', 'tell', 'reduce', 'cut', 'import', 'help', 'ad
d', 'agree', 'offer', 'take']

Most popular Adjectives associate are:  ['japanese', 'foreign', 'major', 'large', 'econom
ic', 'iranian', 'european', 'american', 'new', 'strong']


The top 3 Organizations are:
```

```
In [18]:    1  top_organizations[:3]
```

Out[18]:  [('cts', 6115), ('ec', 883), ('vs', 711)]

And the most frequently occurring terms with these entities are given below.

**Note:** Inside the nested list, the first list is for NOUNS, the second list is for VERBS and third list is for ADJECTIVES

```
In [19]:  1  for entity_name, count in top_organizations[:3]:
          2      print("------------------------------------------------")
          3      print("The entity name is:", entity_name)
          4      print()
          5      print(
          6          "Most popular Nouns associate are: ",
          7          organization_most_popular_terms[entity_name][0],
          8      )
          9      print()
         10      print(
         11          "Most popular Verbs associate are: ",
         12          organization_most_popular_terms[entity_name][1],
         13      )
         14      print()
         15      print(
         16          "Most popular Adjectives associate are: ",
         17          organization_most_popular_terms[entity_name][2],
         18      )
         19      print()
```

```
------------------------------------------------
The entity name is: cts

Most popular Nouns associate are:  ['shr', 'loss', 'mln', 'profit', 'share', 'qtr', 'ne
t', 'dlr', 'div', 'oper']

Most popular Verbs associate are:  ['say', 'set', 'include', 'exclude', 'raise', 'dilut
e', 'end', 'rev', 'declare', 'report']

Most popular Adjectives associate are:  ['net', 'quarterly', '4th', 'prior', 'extraordina
ry', '1st', 'payable', 'regular', '3rd', 'fiscal']

------------------------------------------------
The entity name is: ec

Most popular Nouns associate are:  ['trade', 'u.s', 'community', 'tonne', 'export', 'mini
ster', 'oil', 'tax', 'sugar', 'country']

Most popular Verbs associate are:  ['say', 'propose', 'tell', 'agree', 'meet', 'import',
'add', 'sell', 'offer', 'impose']

Most popular Adjectives associate are:  ['european', 'japanese', 'new', 'foreign', 'frenc
h', 'free', 'non', 'agricultural', 'spanish', 'white']

------------------------------------------------
The entity name is: vs

Most popular Nouns associate are:  ['loss', 'profit', 'ct', 'shr', 'oper', 'mln', 'dlrs',
'net', 'year', 'sale']

Most popular Verbs associate are:  ['include', 'give', 'nil', 'dilute', 'note', 'gain',
'exclude', 'correct', 'discontinue', 'plow']

Most popular Adjectives associate are:  ['net', '4th', 'extraordinary', '1st', '3rd', 'sh
rs', '2nd', 'discontinued', 'public', 'compact']
```

## 5. Extra Credit

# To determine which persons, organizations, and locations most frequently occur in the same sentences.

## a. Approach for solving Extra Credit:

We store the most frequently occuring person, location and organizations in a sentence using the following data structure:

```
{"sentence_1":[
                {person1 : count_of_person1_in_sentence_1, person2: ....},
                {location1 : count_of_location1_in_sentence_1, location2: ....},
                {organization1 : count_of_organization1_in_sentence_1, organization
2: ....},
                ],
  "sentence_2:".....}
```

## Code:

```python
def frequency_of_entities_in_sentence(doc_id, doc_text):

    sentence_dict = {}

    corpus_cleaned = preprocess_doc_text(doc_text)

    # create doc object
    analyzed_doc = nlp(corpus_cleaned)

    for entity in analyzed_doc.ents:
        entity_text = entity.text.strip()
        sentence = entity.sent.text.strip()

        if entity_text != "":
            if sentence not in sentence_dict.keys():
                # not present then initialize
                sentence_dict[sentence] = [
                    {},  # person
                    {},  # location
                    {},  # organization
                ]

            # for persons
            if entity.label_ == "PERSON":
                if entity_text in sentence_dict[sentence][0].keys():
                    # present then increment count of that entity
                    sentence_dict[sentence][0][entity_text] += 1
                else:
                    # not present
                    sentence_dict[sentence][0][entity_text] = 1

            # for locations
            if entity.label_ == "GPE":
                if entity_text in sentence_dict[sentence][1].keys():
                    # present then increment count of that entity
                    sentence_dict[sentence][1][entity_text] += 1
                else:
                    # not present
                    sentence_dict[sentence][1][entity_text] = 1

            # for organizations
            if entity.label_ == "ORG":
                if entity_text in sentence_dict[sentence][2].keys():
                    # present then increment count of that entity
                    sentence_dict[sentence][2][entity_text] += 1
                else:
                    # not present
                    sentence_dict[sentence][2][entity_text] = 1

    return sentence_dict
```

```
In [21]:   1  # Determine which persons, organizations, and locations most frequently occur in the s
           2  combined_sentence_dict = {}
           3
           4  # here we have restricted the code to run for 10 documents only
           5  for doc_id in reuters.fileids()[:10]:
           6
           7      sentence_dicts = frequency_of_entities_in_sentence(
           8          doc_id, reuters.open(doc_id).read()
           9      )
          10
          11      for sentence_with_frequent_entity in sentence_dicts.keys():
          12          if sentence_with_frequent_entity in combined_sentence_dict.keys():
          13              # present
          14              # update the person dictionaries of that sentence
          15              combined_sentence_dict[sentence_with_frequent_entity][0].update(
          16                  sentence_dicts[sentence_with_frequent_entity][0]
          17              )
          18              # update the location dictionaries of that sentence
          19              combined_sentence_dict[sentence_with_frequent_entity][1].update(
          20                  sentence_dicts[sentence_with_frequent_entity][1]
          21              )
          22              # update the organization dictionaries of that sentence
          23              combined_sentence_dict[sentence_with_frequent_entity][2].update(
          24                  sentence_dicts[sentence_with_frequent_entity][2]
          25              )
          26          else:
          27              # not present
          28              # add the combined dictionary
          29              combined_sentence_dict[sentence_with_frequent_entity] = sentence_dicts[
          30                  sentence_with_frequent_entity
          31              ]
```

```
In [22]:   1  def most_frequently_occuring_entities(combined_sentence_dict):
           2      result = {}
           3      for key, value in combined_sentence_dict.items():
           4
           5          entities = [{}, {}, {}]
           6
           7          # person
           8          entities[0] = sorted(value[0].items(), key=lambda item: -item[1])[:5]
           9
          10          # location
          11          entities[1] = sorted(value[1].items(), key=lambda item: -item[1])[:5]
          12
          13          # organization
          14          entities[2] = sorted(value[2].items(), key=lambda item: -item[1])[:5]
          15
          16          result[key] = entities
          17
          18      return result
```

```
In [23]:   1  most_frequently_occuring_entities_in_sentences = most_frequently_occuring_entities(
           2      combined_sentence_dict
           3  )
```

For the first 10 sentence the top 3 most frequently occuring person, location and organization are as follows:

```
In [24]:   1   i = 0
           2   for sentence, entities in most_frequently_occuring_entities_in_sentences.items():
           3       print("--------------------")
           4       print("Sentence is :", sentence)
           5       print()
           6       print("Frequently Occuring Persons in sentence:", entities[0])
           7       print("Frequently Occuring Locations in sentence:", entities[1])
           8       print("Frequently Occuring Organizations in sentence:", entities[2])
           9       print()
          10       # to print only the first 10 sentences
          11       i += 1
          12       if i == 10:
          13           break
```

```
--------------------
Sentence is : ASIAN EXPORTERS FEAR DAMAGE FROM U.S. JAPAN RIFT Mounting trade friction
between the U.S.

Frequently Occuring Persons in sentence: []
Frequently Occuring Locations in sentence: [('U.S.', 2), ('JAPAN', 1)]
Frequently Occuring Organizations in sentence: [('RIFT Mounting', 1)]

--------------------
Sentence is : And Japan has raised fears among many of Asia's exporting nations that t
he row could inflict far reaching economic damage, businessmen and officials said.

Frequently Occuring Persons in sentence: []
Frequently Occuring Locations in sentence: [('Japan', 1)]
Frequently Occuring Organizations in sentence: []

--------------------
Sentence is : They told Reuter correspondents in Asian capitals a U.S. Move against Ja
pan might boost protectionist sentiment in the U.S.

Frequently Occuring Persons in sentence: [('Reuter', 1)]
Frequently Occuring Locations in sentence: [('U.S.', 2), ('Japan', 1)]
Frequently Occuring Organizations in sentence: []

--------------------
Sentence is : And lead to curbs on American imports of their products.

Frequently Occuring Persons in sentence: []
Frequently Occuring Locations in sentence: []
Frequently Occuring Organizations in sentence: []

--------------------
Sentence is : But some exporters said that while the conflict would hurt them in the l
ong run, in the short term Tokyo's loss might be their gain.

Frequently Occuring Persons in sentence: []
Frequently Occuring Locations in sentence: [('Tokyo', 1)]
Frequently Occuring Organizations in sentence: []

--------------------
Sentence is : The U.S. Has said it will impose 300 mln dlrs of tariffs on imports of J
apanese electronics goods on April 17, in retaliation for Japan's alleged failure to s
tick to a pact not to sell semiconductors on world markets at below cost.

Frequently Occuring Persons in sentence: []
Frequently Occuring Locations in sentence: [('U.S.', 1), ('Japan', 1)]
```

```
                  Frequently Occuring Organizations in sentence: []


                  ----------------------
                  Sentence is : Unofficial Japanese estimates put the impact of the tariffs at 10 billio
                  n dlrs and spokesmen for major electronics firms said they would virtually halt export
                  s of products hit by the new taxes.

                  Frequently Occuring Persons in sentence: []
                  Frequently Occuring Locations in sentence: []
                  Frequently Occuring Organizations in sentence: []


                  ----------------------
                  Sentence is : We wouldn't be able to do business,  said a spokesman for leading Japane
                  se electronics firm Matsushita Electric Industrial Co Ltd  .

                  Frequently Occuring Persons in sentence: []
                  Frequently Occuring Locations in sentence: []
                  Frequently Occuring Organizations in sentence: [('Matsushita Electric Industrial Co Lt
                  d', 1)]


                  ----------------------
                  Sentence is : A senior official of South Korea's trade promotion association said the
                  trade dispute between the U.S.

                  Frequently Occuring Persons in sentence: []
                  Frequently Occuring Locations in sentence: [("South Korea's", 1), ('U.S.', 1)]
                  Frequently Occuring Organizations in sentence: [('trade promotion association', 1)]


                  ----------------------
                  Sentence is : And Japan might also lead to pressure on South Korea, whose chief export
                  s are similar to those of Japan.

                  Frequently Occuring Persons in sentence: []
                  Frequently Occuring Locations in sentence: [('Japan', 2), ('South Korea', 1)]
                  Frequently Occuring Organizations in sentence: []
```

## 6. Future Scope

In the future we would like to enhance the code by doing the following tasks:

```
    1. Improve runtime of code
    2. Handle different forms of same entity
    3. Overcome same entities being tagged differently
```

In [ ]:  1