

Hypernym Relationship Extraction

In this example, we will use NLTK and Hearst Pattern for hypernym relationship extraction.

- Firstly, install python environment
- Install NLTK: `pip install nltk`
- Download data distribution for NLTK. Install using NLTK downloader: `nltk.download()` . If cannot download using `nltk.download()` , try download manually from https://github.com/nltk/nltk_data/tree/gh-pages! (attachment:image.png) ([https://github.com/nltk/nltk_data/tree/gh-pages!%5Bimage.png%5D\(attachment:image.png\)](https://github.com/nltk/nltk_data/tree/gh-pages!%5Bimage.png%5D(attachment:image.png))) or https://pan.baidu.com/s/1wONWpaa86_wnslksKda8eQ (https://pan.baidu.com/s/1wONWpaa86_wnslksKda8eQ) (code:tfon)
- Unzip the downloaded file to the following folder: `nltk.data.find(". ")`
- Unzip each zip file in the ten folders: *chunkers, corpora, grammars, help, misc, models, sentiment, stemmers, taggers, tokenizers*

Hyponym Extraction using Hearst Pattern

Hyponym extraction follows the following 4 steps:

- Noun phrase chunking or named entity chunking. You can use any np chunking/named entity technique.
- Chunked sentences prepare. Traverse the chunked result, if the label is NP , then merge all the words in this chunk and add a prefix NP_ (for subsequence process).
- Chunking refinement. If two or more NPs next to each other should be merged into a single NP. Eg., "*NP_foo NP_bar blah blah*" becomes "*NP_foo_bar blah blah*"
- Find the hypernym and hyponym pairs based on the refined prepared chunked sentence.

```
In [1]: 1 import nltk
        2 import re
        3 from nltk import pos_tag, word_tokenize, Tree
        4 from nltk.stem import WordNetLemmatizer
```

Regular expression practice: In this example, we show one regex pattern example for Hearst pattern: NP such as `{NP,}* {(or | and)} NP` (<https://docs.python.org/3/library/re.html>) (<https://docs.python.org/3/library/re.html>)

In [2]:

```
1 regex = r"(NP_\w+ (, )?such as (NP_\w+ ?(, )?(and |or )?)+)"
2 test_str = "NP_1 such as NP_2 , NP_3 and NP_4 "
3 matches = re.search(regex, test_str)
4 if matches:
5     # Match.group([group1, ...]) Returns one or more subgroups of the match.
6     # If there is a single argument, the result is a single string;
7     # if there are multiple arguments, the result is a tuple with one item per argument.
8     # Without arguments, group1 defaults to zero (the whole match is returned).
9     print(matches.group(0))
```

NP_1 such as NP_2 , NP_3 and NP_4

Step1: Chunking Sentence

- Note the result is not the chunked np, instead is the chunk tree structure

```
In [3]: 1 def np_chunking(sentence):
2         # your implementation
3
4         return result
5
6 print(np_chunking("""I like to listen to music from musical genres,such as blues,rock and jazz."""))
```

```
(S
  I/PRP
  like/VBP
  to/TO
  listen/VB
  to/TO
  (NP music/NN)
  from/IN
  (NP musical/JJ genres/NNS)
  ,/,
  such/JJ
  as/IN
  (NP blues/NNS)
  ,/,
  (NP rock/NN)
  and/CC
  (NP jazz/NN)
  ./.)
```

Step2: Prepare the chunked result for subsequent Hearst pattern matching

- Traverse the chunked result, if the label is NP , then merge all the words in this chunk and add a prefix NP_
- All the tokens are separated with a white space (" ")
- Remember to lemmatize words, using WordNetLemmatizer (from nltk.stem import WordNetLemmatizer)

```
In [4]: 1 # prepare the chunked sentence by merging words and add prefix NP_
2 def prepare_chunks(chunks):
3     # If chunk is NP, start with NP_ and join tokens in chunk with _ ; Else just keep the token as it is
4     terms = []
5     for chunk in chunks:
6         label = None
7         try:
8             # see if the chunk is simply a word or a NP. But non-NP fail on this method call
9             label = chunk.label()
10        except:
11            pass
12        # Based on the label, do processing, your implementation here...
```

```
In [5]: 1 raw_text = "I like to listen to music from musical genres,such as blues,rock and jazz."
2 chunk_res = np_chunking(raw_text)
3 print(prepare_chunks(chunk_res))
```

I like to listen to NP_music from NP_musical_genre , such as NP_blue , NP_rock and NP_jazz .

Step3: Refinement chunking

If two or more NPs next to each other should be merged into a single NP. E.g., NP_foo NP_bar blah blah becomes NP_foo_bar blah blah

```
In [6]: 1 def merge_NP(prepared_chunks):
2     sentence = re.sub(r"(NP_\w+ NP_\w+)+",lambda m: m.expand(r'\1').replace(" NP_", "_"),prepared_chunks)
3     return sentence
```

```
In [7]: 1 merge_NP("NP_foo NP_bar blah blah")
```

```
Out[7]: 'NP_foo_bar blah blah'
```

Step4: Find the hypernym and hyponyms on processed chunked results

- Define Hearst patterns. Besides the regex, we also need to specify whether the hypernym is in the first part or the second part in the pattern.
 - For example, in the pattern NP1 such as NP2 AND NP3 , the hypernym is the first part of the pattern; in the pattern NP1 , NP2 and other NP3 , the hypernym is the last part of the pattern.

- After regex matching, find all the NPs and extract the hypernym and hyponym pairs based on the first or last attribute.
- Clean the NPs by removing the prefix NP_ and _

In [8]:

```
1 # Given by the prepared text, return the hypernym-hyponym pairs
2 def hyponym_extract(prepared_text, hearst_patterns):
3     # your implementation
4     return pairs
5
6 hearst_patterns = [("(NP_\w+ (, )?such as (NP_\w+ ?(, )?(and |or )?)+"), "first"),
7                    ("((NP_\w+ ?(, )?)+(and |or )?other NP_\w+)", "last")] # two examples for hearst pattern
8 print(hyponym_extract(prepare_chunks(np_chunking("I like to listen to music from musical genres,such as blues,rock and
9 print(hyponym_extract(prepare_chunks(np_chunking("He likes to play basketball,football and other sports.")),hearst_patterns))
```

```
[('NP_blue', 'NP_musical_genre'), ('NP_rock', 'NP_musical_genre'), ('NP_jazz', 'NP_musical_genre')]
[('NP_basketball', 'NP__sport'), ('NP_football', 'NP__sport')]
```

In [10]:

```
1 def find_hyponyms(sentence, hearst_patterns):
2     # your implementation
3
4 print(find_hyponyms("""I like to listen to music from musical genres,such as blues,rock and jazz."""), hearst_patterns)
5 print(find_hyponyms("""He likes to play basketball,football and other sports."""),hearst_patterns))
```

```
[('NP_blue', 'NP_musical_genre'), ('NP_rock', 'NP_musical_genre'), ('NP_jazz', 'NP_musical_genre')]
[('NP_basketball', 'NP__sport'), ('NP_football', 'NP__sport')]
```

In [11]:

```
1 def clean_np(term):
2     return term.replace("NP_", "").replace("_", " ")
3 clean_np('NP_football')
```

Out[11]: 'football'

Complete Program for Hypernym extraction using Hearst Pattern

```
In [12]: 1 # Merge everything to get the final extractor
2 class HearstPatterns(object):
3     # finish the extractor class using the aforementioned functions
4
```

```
In [13]: 1 # Test case for hearst patterns
2 hp = HearstPatterns()
3 test = ["Agar is a substance prepared from a mixture of red algae, such as Gelidium,for laboratory or industrial use.",
4         "... works by such authors as Herrick, Goldsmith, and Shakespeare.",
5         "... bistros, coffee shops, and other cheap eating places.",
6         "...all common law countries, including Canada and England.",
7         "...most European countries, especially France, England, and Spain."]
8         # text = 'I like to listen to music from musical genres such as blues, rock and jazz. He likes to play
9 for txt in test:
10     hps = hp.find_hyponyms(txt)
11     print(hps)
```

```
[('Gelidium', 'red algae')]
[('Herrick', ' author'), ('Goldsmith', ' author'), ('Shakespeare', ' author')]
[('bistro', ' cheap eating place'), ('coffee shop', ' cheap eating place')]
[('Canada', 'common law country'), ('England', 'common law country')]
[('France', 'European country'), ('England', 'European country'), ('Spain', 'European country')]
```