

# Xử lý ngôn ngữ tự nhiên

## Chương 4: Thư viện NLTK, Wordnet

Khoa CNTT, Đại học Kỹ thuật - Công nghệ Cần Thơ  
Lưu hành nội bộ

# Các tác vụ xử lý text cơ bản

- 1 Tokenizing văn bản thành câu
- 2 Tokenizing văn bản thành chữ
- 3 Loại bỏ stopwords trong câu đã được tokenized
- 4 Tìm synsets của 1 từ trong WordNet
- 5 Tìm kiếm ngữ nguyên (lemma), đồng nghĩa (synonyms), trái nghĩa (antonyms)
- 6 Tính toán sự giống nhau của từ
- 7 Khám phá sự kết hợp (collocations)
- 8 Stemming, lemmatization, Part-of-speech

# Tìm synsets của một từ trong WordNet

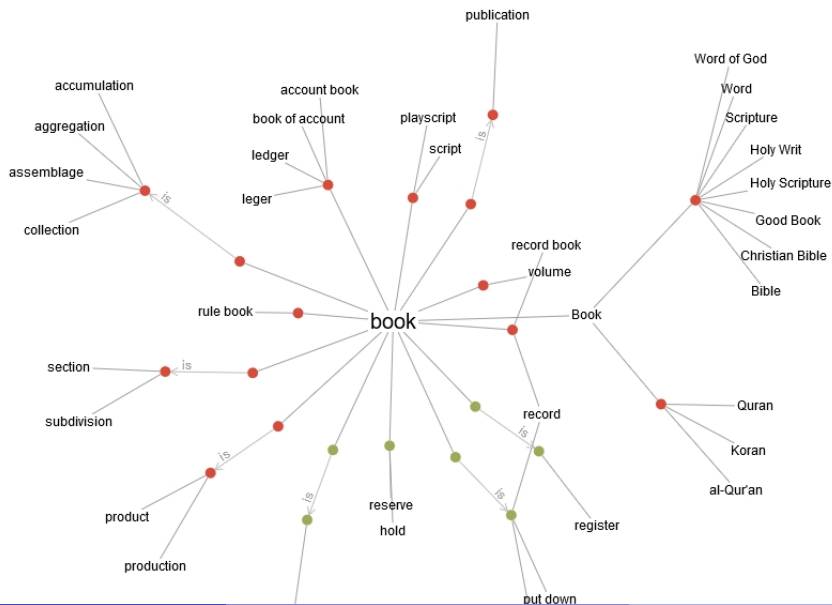
- WordNet là cơ sở dữ liệu từ vựng học của ngôn ngữ tiếng anh
- WordNet là từ điển thiết kế cụ thể cho NLP
- WordNet cung cấp giao diện đơn giản để tìm kiếm, tra từ
- Kết quả trả về là một danh sách (list) các Synset
- Synsets được sắp xếp theo cấu trúc tương tự như một cây thừa kế (inheritance tree)

# Tìm synsets của một từ trong WordNet

```
1 from nltk.corpus import wordnet
2 syn = wordnet.synsets('cookbook')[0]
3 print(syn.name())
4 print(syn.definition())
5
6 #cookbook.n.01
7 #a book of recipes and cooking directions
```

```
1 from nltk.corpus import wordnet
2 syns = wordnet.synsets('book')
3 print(len(syns))
4 for syn in syns:
5     print(syn.name())
6     print(syn.definition())
7
8 #15
9 #book.n.01
10 #a written work or composition that has been published (printed on
    pages bound together)
11 #book.n.02
12 #physical objects consisting of a number of pages bound together
13 #record.n.05
14 #a compilation of the known facts regarding something or someone
15 #script.n.01
16 #a written version of a play or other dramatic composition; used in
    preparing for a performance
17 # ...
```

```
1 from nltk.corpus import wordnet
2 syns = wordnet.synsets('book')
3 print(len(syns))
4 for syn in syns:
5     print(syn.examples())
6
7 #15
8 #['I am reading a good book on economics']
9 #
10 #['he used a large book as a doorstop']
11 #
12 #["Al Smith used to say, 'Let's look at the record'", 'his name is in
    all the record books']
13 #
14 #[]
15 #...
```



# Hypernyms

- Do synsets được cấu trúc dạng cây thừa kế
  - nút cha được gọi là hypernyms: trừu tượng
  - nút con được gọi là hyponyms: cụ thể
  - các nút con có ý nghĩa gần tương tự nhau
- Synsets có thể truy ngược lên nút root



# Hypernyms, hyponyms

```
1 from nltk.corpus import wordnet
2 syn = wordnet.synsets('cookbook')[0]
3
4 print(syn.hypernyms())
5 print()
6 print(syn.hypernyms()[0].hyponyms())
7 print()
8 print(syn.root_hypernyms())
9 print()
10 print(syn.hypernym_paths())
11
12 [Synset('reference_book.n.01')]
```

13

```
14 [Synset('annual.n.02'), Synset('atlas.n.02'), Synset('cookbook.n.01'),  
    Synset('directory.n.01'), Synset('encyclopedia.n.01'), Synset('handbook.n.01'),  
    Synset('instruction_book.n.01'), Synset('source_book.n.01'), Synset('wordbook.n.01')]
```

15

```
16 [Synset('entity.n.01')]
```

17

```
18 [[Synset('entity.n.01'), Synset('physical_entity.n.01'), Synset('object.n.01'),  
    Synset('whole.n.02'), Synset('artifact.n.01'), Synset('creation.n.02'),  
    Synset('product.n.02'), Synset('work.n.02'), Synset('publication.n.01'),  
    Synset('book.n.01'), Synset('reference_book.n.01'), Synset('cookbook.n.01')]]
```

# Part of speech (POS)

- Một synsets của một từ là một list
- Trong list trả về trên, có nhiều từ loại như danh từ (n), tính từ (a), trạng từ (r), động từ (v)

# Part of speech (POS)

```
1 from nltk.corpus import wordnet
2 syn = wordnet.synsets('great')[0]
3
4 print(syn.pos())
5 print(len(wordnet.synsets('great')))
6 print(len(wordnet.synsets('great', pos='n')))
7 print(len(wordnet.synsets('great', pos='a')))
8
9 #n
10 #7
11 #1
12 #6
```

# Lemmas

- Lemma: dạng chuẩn tắc (canonical form) hoặc dạng ngữ nguyên (morphological form) của từ (word)

```
1 from nltk.corpus import wordnet
2 syn = wordnet.synsets('cookbook')[0]
3
4 lemmas = syn.lemmas()
5 print(len(lemmas))
6 print(lemmas[0].name())
7 print(lemmas[1].name())
8 print(lemmas[0].synset() == lemmas[1].synset())
9
10 lems = [lemma.name() for lemma in syn.lemmas()]
11 print(lems)
12
13 #2
14 #cookbook
15 #cookery_book
16 #True
17 #['cookbook', 'cookery_book']
```

# Lemmas

- Mỗi Synsets là một danh sách (list), mỗi element trong list có một hoặc nhiều lemmas
- Tìm tất cả các lemmas của Synsets
- Sẽ có nhiều lemmas trùng nhau, loại bỏ bằng cách đổi kiểu list thành kiểu tập hợp (set)

# Lemmas

```
1 from nltk.corpus import wordnet
2
3 synonyms = []
4 for syn in wordnet.synsets('book'):
5     for lemma in syn.lemmas():
6         synonyms.append(lemma.name())
7
8 print(len(synonyms))
9 print(len(set(synonyms)))
10 print(synonyms)
11 #38
12 #25
13 #['book', 'book', 'volume', 'record', 'record_book', 'book', 'script',
    'book', 'playscript', 'ledger', 'leger', #'account_book', '
    book_of_account', 'book', 'book', 'book', 'rule_book', 'Koran', '
    Quran', "al-Qur'an", 'Book', #'Bible', 'Christian_Bible', 'Book', '
    ...]
```



# Autonyms - từ khác nghĩa

```
1 from nltk.corpus import wordnet
2
3 antonyms = []
4
5 for syn in wordnet.synsets("small"):
6     for l in syn.lemmas():
7         if l.antonyms():
8             antonyms.append(l.antonyms()[0].name())
9
10 print(antonyms)
11 print(set(antonyms))
12
13 #['large ', 'big ', 'big ']
14 #{'big ', 'large '}
```

# synonyms - từ đồng nghĩa

```
1 from nltk.corpus import wordnet
2
3 synonyms = []
4
5 for syn in wordnet.synsets('small'):
6     for lemma in syn.lemmas():
7         synonyms.append(lemma.name())
8
9 print(set(synonyms))
10
11 #{'modest', 'pocket-size', 'diminished', 'small-scale', 'lowly', '
    pocket-sized', 'minor', 'low', 'belittled', #'little', 'small', '
    humble', 'minuscule'}
```

# Tính sự giống nhau của từ

- Do Synsets có cấu trúc cây, nên 2 synsets càng gần nhau thì càng giống nhau (similarity) về ý nghĩa

```
1 from nltk.corpus import wordnet
2 cb = wordnet.synsets('cookbook')[0]
3 ib = wordnet.synsets('instruction_book')[0]
4 print(cb.wup_similarity(ib))
5
6 #0.9166666666666666
7
8 table = wordnet.synsets('table')[0]
9 print(table.wup_similarity(ib))
10
11 #0.1111111111111111
```

- Synsets table và synsets instruction\_book giống nhau 11%
- Kiểm tra các hypernyms chung của 2 synsets

```
12 print(sorted(table.common_hypernyms(ib)))  
13  
14 #[Synset('entity.n.01')]
```

# Tìm collocations

- Collocations là 2 hoặc nhiều word xuất hiện gần nhau thường xuyên, ví dụ word Cần và word Thơ
- Tại ngữ cảnh khác nhau thì collocations cũng khác nhau
- Thư viện nltk hỗ trợ: BigramCollocationFinder và TrigramCollocationFinder

```

1 from nltk.corpus import webtext
2 from nltk.collocations import BigramCollocationFinder
3 from nltk.metrics import BigramAssocMeasures
4 from nltk.corpus import stopwords
5
6 words = [w.lower() for w in webtext.words('singles.txt')]
7
8 stopset = set(stopwords.words('english'))
9 filter_stops = lambda w: len(w) < 3 or w in stopset
10
11 bcf = BigramCollocationFinder.from_words(words)
12
13 bcf.apply_word_filter(filter_stops)
14 # appear more than 4 times
15 print(bcf.nbest(BigramAssocMeasures.likelihood_ratio, 4))
16
17 #[( 'would', 'like'), ('age', 'open'), ('medium', 'build'), ('social', '
    drinker')]

```

# Stemming words

- Kỹ thuật gỡ bỏ đuôi của từ (ed, ing)
  - PorterStemmer
  - LancasterStemmer
  - SnowballStemmer
  - RegexpStemmer

```
1 from nltk.stem import PorterStemmer
2 stemmer = PorterStemmer()
3 print(stemmer.stem('cooking'))
4 print(stemmer.stem('cookery'))
5
6 #cook
7 #cookeri
8
9 from nltk.stem import LancasterStemmer
10 stemmer = LancasterStemmer()
11 print(stemmer.stem('cooking'))
12 print(stemmer.stem('cookery'))
13
14 #cook
15 #cookery
```



```
1 from nltk.stem import SnowballStemmer
2 stemmer = LancasterStemmer()
3 print(stemmer.stem('cooking'))
4 print(stemmer.stem('cookery'))
5
6 #cook
7 #cookery
8
9 from nltk.stem import RegexpStemmer
10 stemmer = RegexpStemmer('ing')
11 print(stemmer.stem('cooking'))
12 print(stemmer.stem('cookery'))
13 print(stemmer.stem('ingleside'))
14
15 #cook
16 #cookery
17 #leside
```

# Khác nhau giữa Stemming và Lemmatization

```
1 from nltk.stem import WordNetLemmatizer
2 from nltk.stem import PorterStemmer
3
4 stemmer = PorterStemmer()
5 lemmatizer = WordNetLemmatizer()
6
7 print(stemmer.stem('stones'))
8 print(stemmer.stem('speaking'))
9 print(stemmer.stem('bedroom'))
10 print(stemmer.stem('jokes'))
11 print(stemmer.stem('lisa'))
12 print(stemmer.stem('purple'))
13 print('_____')
```

```
14  
15 print(lemmatizer.lemmatize('stones'))  
16 print(lemmatizer.lemmatize('speaking'))  
17 print(lemmatizer.lemmatize('bedroom'))  
18 print(lemmatizer.lemmatize('jokes'))  
19 print(lemmatizer.lemmatize('lisa'))  
20 print(lemmatizer.lemmatize('purple'))
```

21

22 stone

23 speak

24 bedroom

25 joke

26 lisa

27 purpl

28

---

29 stone

30 speaking

31 bedroom

32 joke

33 lisa

34 purple

# Thay thế word bằng regular expressions

- Tạo file replacers.py có nội dung như sau:

```
1 import re
2
3 replacement_patterns = [
4     (r'won\'t', 'will not'),
5     (r'can\'t', 'cannot'),
6     (r'i\'m', 'i am'),
7     (r'ain\'t', 'is not'),
8     (r'(\w+)\'ll', '\g<1> will'),
9     (r'(\w+)\n\'t', '\g<1> not'),
10    (r'(\w+)\\'ve', '\g<1> have'),
11    (r'(\w+)\'s', '\g<1> is'),
12    (r'(\w+)\\'re', '\g<1> are'),
13    (r'(\w+)\'d', '\g<1> would'),
14 ]
```

```
15
16 class RegexpReplacer(object):
17     def __init__(self, patterns=replacement_patterns):
18         self.patterns = [(re.compile(regex), repl) for (regex, repl) in
19                             patterns]
20
21     def replace(self, text):
22         s = text
23
24         for (pattern, repl) in self.patterns:
25             s = re.sub(pattern, repl, s)
26
27     return s
```

- File chương trình chính cùng cấp với file replacers.py

```
1 from replacers import RegexpReplacer
2 replacer = RegexpReplacer()
3 print(replacer.replace("can't is a contraction"))
4 print(replacer.replace("I should've done that thing I didn't do"))
5
6 #cannot is a contraction
7 #I should have done that thing I did not do
```

- Thay thế ký tự lặp lại (repeating characters): sinh viên xem bài cũ

# Data kèm theo thư viện NLTK

- Xem lại bài trước về cách cài đặt data của thư viện NLTK
- Ví dụ data được cài vào thư mục `D:\nltk_data`
- Thư mục `D:\nltk_data\corpora` chứa các data sau khi cài đặt



```
1 from nltk.corpus import names
2 print(names.fileids())
3 print(len(names.words('female.txt')))
4 print(len(names.words('male.txt')))
5
6 from nltk.corpus import words
7 print(words.fileids())
8 print(len(words.words('en-basic')))
9 print(len(words.words('en')))
10
11 #['female.txt', 'male.txt']
12 #5001
13 #2943
14 #['en', 'en-basic']
15 #850
16 #235886
```

# Wordlist do người dùng tự tạo

- Tạo thư mục `nlp-corpus` trong `D:\nltk_data\corpora`
- Tạo file `mywords.txt`, lưu trong thư mục `nlp-corpus`, với nội dung sau:

```
1 nltk
2 corpus
3 corpora
4 wordnet
```

```
1 import nltk.data
2 from nltk.corpus.reader import WordListCorpusReader
3
4 reader = WordListCorpusReader('D:/nltk_data/corpora/nlp-corpus/', [
    'mywords.txt'])
5 print(reader.words())
6 print(reader.fileids())
7
8 #['nltk', 'corpus', 'corpora', 'wordnet']
9 #['mywords.txt']
```

# Part-of-speech tagging

- Xem lại bài cũng đôi với tác vụ này
- Trong ví dụ này, chúng ta sẽ load từ file đã có nội dung word/tag
- Tạo file `txt-pos.txt` trong thư mục  
`D:\nltk_data\corpora\nlp-corpus` đã tạo ở slides trước
- Nội dung file: `The/at-tl expense/nn and/cc time/nn  
involved/vbn are/ber astronomical/jj ./.`

```
1 from nltk.corpus.reader import TaggedCorpusReader
2 reader = TaggedCorpusReader("D:/nltk_data/corpora/nlp-corpus", ["txt-
    pos.txt"], tagset="en-brown")
3 print(reader.words())
4 print(reader.sents())
5 print(reader.paras())
6 print(reader.tagged_words(tagset="universal"))
7 print(reader.tagged_sents(tagset="universal"))
8 print(reader.tagged_paras(tagset="universal"))
```

```
9
10 #['The', 'expense', 'and', 'time', 'involved', 'are', ...]
11 #[['The', 'expense', 'and', 'time', 'involved', 'are', 'astronomical']]
12 #[[['The', 'expense', 'and', 'time', 'involved', 'are', 'astronomical'
    '']]]
13 #[( 'The', None), ('expense', None), ('and', None), ...]
14 #[( 'The', None), ('expense', None), ('and', None), ('time', None), ('
    involved', None), ('are', None), #('astronomical', None)]]
15 #[[[( 'The', None), ('expense', None), ('and', None), ('time', None), ('
    involved', None), ('are', None), #('astronomical', None)]]]
```

# Sử dụng WordNet cho tác vụ tagging

- Tạo file taggers.py với nội dung như sau:

```
1 from nltk.tag import NgramTagger, SequentialBackoffTagger
2 from nltk.corpus import wordnet
3 from nltk.probability import FreqDist
4
5 class WordNetTagger(SequentialBackoffTagger):
6     def __init__(self, *args, **kwargs):
7         SequentialBackoffTagger.__init__(self, *args, **kwargs)
8         self.wordnet_tag_map = {
9             'n': 'NN',
10            's': 'JJ',
11            'a': 'JJ',
12            'r': 'RB',
13            'v': 'VB'
14        }
```

```
15 def choose_tag(self, tokens, index, history):
16     word = tokens[index]
17     fd = FreqDist()
18
19     for synset in wordnet.synsets(word):
20         fd[synset.pos()] += 1
21
22     if not fd: return None
23     return self.wordnet_tag_map.get(fd.max())
```



- Tạo file chương trình chính cùng cấp với file taggers.py

```
1 from taggers import WordNetTagger
2 wt = WordNetTagger()
3 print(wt.tag(['food', 'is', 'great']))
4
5 #[( 'food ', 'NN'), ( 'is ', 'VB'), ( 'great ', 'JJ')]
```