

# Handout 5: Corpora and Lexica

## Corpora

### 1. Corpus objects

- a. Reside in module `nltk.corpus`

```
1 >>> from nltk.corpus import gutenberg
2 >>> from nltk.corpus import brown
```

- b. Corpus description: `print(gutenberg.readme())`

### 2. Characters, words, sentences

```
1 >>> gutenberg.raw()
2 '[Emma by Jane Austen ...'
3 >>> gutenberg.words()
4 ['[', 'Emma', 'by', 'Jane', 'Austen', '1816', ']', ...]
5 >>> gutenberg.sents()
6 [['[', 'Emma', 'by', 'Jane', 'Austen', '1816', ']', ['VOLUME', 'I'], ...]
```

### 3. How to avoid typing “gutenberg” all the time?

```
1 >>> from nltk.corpus import gutenberg as gb
2 >>> gb.words()
3 ['[', 'Emma', 'by', 'Jane', 'Austen', '1816', ']', ...]
```

### 4. Selecting a single file

```
1 >>> gb.words('melville-moby_dick.txt')
2 ['[', 'Moby', 'Dick', 'by', 'Herman', 'Melville', ...]
3 >>> gb.fileids()
4 ['austen-emma.txt', 'austen-persuasion.txt', 'austen-sense.txt', ...]
```

### 5. Instead of doing `from nltk.book import *`

```
1 >>> from nltk import Text
2 >>> text1 = Text(gb.words('melville-moby_dick.txt'))
```

### 6. *Review.* How do we do the following?

- a. Get the number of tokens.
- b. Get the wordlist in alphabetic order.
- c. Get the word frequency distribution.
- d. Get the relative frequency of word  $w$ .
- e. Get the most-frequent  $n$  words.
- f. Make a concordance of word  $w$ .

## 7. Categories

- a. How do I get a genre of the Brown corpus?

```
1 from nltk.corpus import brown
2 brown.words(categories='news')
```

- b. What if I forget what the genres are?

```
1 >>> brown.categories()
2 ['adventure', 'belles_lettres', 'editorial', ..., 'science_fiction']
```

- c. Making one big text out of multiple selected genres

```
brown.words(categories=['news', 'editorial'])
```

- d. Files in a category: `brown.fileids(categories='romance')`

- e. Categories for a file: `brown.categories(fileids='cp27')`

- f. Not all corpora have categories: `gutenberg.categories()` – error

- g. Reuters corpus has overlapping categories

```
>>> reuters.categories('training/9865')
['barley', 'corn', 'grain', 'wheat']
>>> reuters.fileids(categories='corn')
['test/14832', 'test/14858', 'test/15033', ..., 'training/9989']
```

- h. Temporal information in filename:

```
1 >>> [fid[:4] for fid in inaugural.fileids()]
2 ['1789', '1793', '1797', '1801', '1805', '1809', ..., '2009']
```

## 8. Corpora discussed in the chapter

Corpus	Files	FileIDs	Categories
<b>gutenberg</b>	books	<i>author-titleword.txt</i>	–
<b>webtext</b>	misc from web	<i>titleword.txt</i>	–
<b>brown</b>	representative corpus	<i>c genre no</i>	genres
<b>reuters</b>	newswire articles	<i>test/no, training/no</i>	topics
<b>inaugural</b>	inaugural addresses	<i>year-pres.txt</i>	–
<b>cess_esp</b>	spanish	<i>code.tbf</i>	–
<b>floresta</b>	portuguese	[only one file]	–
<b>indian</b>	one file per language	<i>lang.pos</i>	–
<b>udhr</b>	one file per language	<i>lang-charenc</i>	–

## Lexical resources

### 9. English wordlist

#### a. nltk.corpus.words

```
1 >>> from nltk.corpus import words as wordlist
2 >>> wordlist.words()
3 ['A', 'a', 'aa', 'aal', 'aalii', 'aam', 'Aani', 'aardvark', ...]
```

#### b. Example:

```
1 >>> eng_vocab = set(w.lower() for w in wordlist.words())
```

#### c. Finding words that are not in the dictionary

```
1 >>> [w for w in set(text1) if w.lower() not in eng_vocab]
2 ['woods', 'clotted', 'plaudits', 'marching', 'disobeying', 'canes', ...]
```

#### d. How do we strip the endings? → *stemmers*

### 10. Other wordlists

#### a. Recall the highest-frequency pairs in *Moby Dick*: ('', 'and'), ('of', 'the'), ... We'd like to get rid of the “empty” words.

#### b. Stopwords

```
1 >>> from nltk.corpus import stopwords
2 >>> stopwords.fileids()
3 ['danish', 'dutch', 'english', 'finnish', 'french', 'german', ...]
4 >>> stopwords.words('english')
5 ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', ..., 'now']
```

#### c. Names

```
1 >>> from nltk.corpus import names
2 >>> names.fileids()
3 ['female.txt', 'male.txt']
4 >>> names.words('female.txt')
5 ['Abagael', 'Abigail', 'Abbe', 'Abbey', 'Abbi', 'Abbie', ...]
```

### 11. CMU pronunciation dictionary

#### a. >>> from nltk.corpus import cmudict

```
2 >>> cmudict.words()
3 ['a', 'a.', 'a', 'a42128', 'aaa', 'aaberg', ...]
```

#### b. New method: entries()

```
1 >>> cmudict.entries()[:5]
2 [('a', ['AH0']), ('a.', ['EY1']), ('a', ['EY1']), ('a42128', ['EY1',
3 'F', 'AO1', 'R', 'T', 'UW1', 'W', 'AH1', 'N', 'T', 'UW1', 'EY1',
4 'T']), ('aaa', ['T', 'R', 'IH2', 'P', 'AH0', 'L', 'EY1'])]
```

- c. Each entry is a pair (word, pronunc)

```

1 >>> for (word, pronunc) in cmudict.entries():
2     ...     if word.endswith('tale'):
3         ...         print(word, pronunc)
4     ...
5     artale ['AA0', 'R', 'T', 'AA1', 'L', 'IYO']
6     denatale ['D', 'IH0', 'N', 'AA0', 'T', 'AA1', 'L', 'IYO']
7     ...
8     vitale ['V', 'IH0', 'T', 'AE1', 'L', 'IYO']
9     vitale ['V', 'AY2', 'T', 'AE1', 'L']

```

- d. The numbers represent stress: 1 > 2 > 0

```

1 aaa ['T', 'R', 'IH2', 'P', 'AH0', 'L', 'EY1']
2 fairytale ['F', 'EH1', 'R', 'IYO', 'T', 'EY2', 'L']

```

## 12. Using a dict for look-up

- a. Create from a list of pairs

```

1 >>> pron_dict = dict(cmudict.entries())
2 >>> pron_dict['apple']
3 ['AE1', 'P', 'AH0', 'L']

```

- b. `cmudict.dict()` – values are *lists* of pronunciations

## 13. Example: rhymes with

- a. Words that rhyme with “cat”:

```

1 >>> pron_dict['cat']
2 ['K', 'AE1', 'T']
3 >>> [w for (w,p) in cmudict.entries()
4     ...     if len(p) > 2 and p[-2:] == ['AE1', 'T']]
5 ['arnatt', 'at-bat', 'balyeat', 'bat', 'batt', 'batte', 'begat', ...]

```

- b. The “rhyming suffix” of a word starts with the (last) stressed vowel

```

1 >>> pron = ['T', 'AH1', 'M', 'K', 'AE1', 'T']
2 >>> list(enumerate(pron))
3 [(0, 'T'), (1, 'AH1'), (2, 'M'), (3, 'K'), (4, 'AE1'), (5, 'T')]
4 >>> [i for (i,ph) in enumerate(pron) if '1' in ph]
5 [1,4]
6 >>> _[-1]
7 4
8 >>> i = _
9 >>> suf = pron[i:]
10 >>> suf
11 ['AE1', 'T']

```

- c. Find words that end with the rhyming suffix

```

1 >>> [w for (w,p) in cmudict.entries() if p[-len(suf):] == suf]
2 ['arnatt', 'at', 'at-bat', 'balyeat', 'bat', 'batt', ..., 'vat']

```

## 14. Swadesh lists

- a. 100 common words in multiple languages

```
1 >>> from nltk.corpus import swadesh
2 >>> swadesh.fileids()
3 ['be', 'bg', 'bs', 'ca', 'cs', 'cu', 'de', 'en', 'es', 'fr', 'hr',
4  'it', 'la', 'mk', 'nl', 'pl', 'pt', 'ro', 'ru', 'sk', 'sl', 'sr',
5  'sw', 'uk']
```

- b. Each entry is a tuple corresponding to one meaning

```
1 >>> swadesh.entries()
2 [('\\xd1\\x8f', '\\xd0\\xb0\\xd0\\xb7', 'ja', ..., 'I', ...),
3  ('\\xd1\\x82\\xd1\\x8b', '\\xd1\\x82\\xd0\\xb8', 'ti', ..., 'you (singular), thou', .
4  ...]
```

- c. Can select just some languages

```
1 >>> swadesh.entries(['en', 'de'][:5])
2 [('I', 'ich'), ('you (singular), thou', 'du, Sie'), ('he', 'er'), ...]
```

- d. Make a table for translation

```
1 >>> trans = dict(swadesh.entries(['en', 'de']))
2 >>> trans['dog']
3 'Hund'
```

## Wordnet

### 15. Words, synsets, lemmas

- a. Each word has some number of “meanings,” called **synsets**

```
1 >>> from nltk.corpus import wordnet as wn
2 >>> horse = wn.synsets('horse')
3 >>> horse
4 [Synset('horse.n.01'), Synset('horse.n.02'), Synset('cavalry.n.01'),
5  Synset('sawhorse.n.01'), Synset('knight.n.02'), Synset('horse.v.01')]
```

- b. Meanings represented indirectly by the **set** of **synonyms** that share the meaning

```
1 >>> horse[0].lemma_names()
2 ['horse', 'Equus_caballus']
3 >>> horse[1].lemma_names()
4 ['horse', 'gymnastic_horse']
5 >>> horse[4].lemma_names()
6 ['knight', 'horse']
```

**c.** Synonyms share a synset

```
1 >>> knight = wn.synsets('knight')
2 >>> [s.lemma_names() for s in knight]
3 [['knight'], ['knight', 'horse'], ['knight', 'dub']]
4 >>> knight[1] == horse[4]
5 True
```

**d.** Each synset has a unique part of speech

```
1 >>> horse[0].pos()
2 'n'
3 >>> horse[5].pos()
4 'v'
```

**e.** Restricting synsets to a pos

```
1 >>> wn.synsets('horse', 'v')
2 [Synset('horse.v.01')]
```

**16.** Relations between synsets

**a.** Hypernyms and hyponyms

```
1 >>> horse[0].hypernyms()
2 [Synset('equine.n.01')]
3 >>> horse[0].hyponyms()
4 [Synset('hack.n.07'), Synset('stalking-horse.n.04'),
5  Synset('saddle_horse.n.01'), Synset('pacer.n.02'), ...]
```

**b.** Part meronyms and holonyms

```
1 >>> cell = wn.synsets('cell', 'n')
2 >>> cell[1].part_meronyms()
3 [Synset('cytoplasm.n.01'), Synset('nucleus.n.01'), Synset('organelle.n.01'),
4  Synset('cell_membrane.n.01'), Synset('energid.n.01'), Synset('vacuole.n.01')]
5 >>> _[0].part_holonyms()
6 [Synset('cell.n.02')]
```

**c.** Substance meronyms and holonyms

```
1 >>> animal = wn.synsets('animal', 'n')
2 >>> animal[0].substance_meronyms()
3 [Synset('animal_tissue.n.01')]
```

**d.** Member meronyms and holonyms

```
1 >>> oscines = wn.synsets('oscines', 'n')
2 >>> oscines[0].member_meronyms()
3 [Synset('xenicidae.n.01'), Synset('hirundinidae.n.01'), Synset('oscine.n.01'),
```

**e.** Entailments

```
1 >>> snore = wn.synsets('snore', 'v')
2 >>> snore[0].entailments()
3 [Synset('sleep.v.01')]
```

f. Attribute

```
1 >>> fast = wn.synsets('fast', 'a')
2 >>> fast[0].attributes()
3 [Synset('speed.n.02')]
```

17. Lemmas

a. A lemma is the pairing of a word and a synset: it is a **word sense**

```
1 >>> horse_senses = wn.lemmas('horse')
2 >>> horse_senses[4]
3 Lemma('knight.n.02.horse')
4 >>> horse_senses[4].name()
5 'horse'
6 >>> horse_senses[4].synset()
7 Synset('knight.n.02')
```

b. From synset to lemmas expressing it

```
1 >>> horse[4].lemmas()
2 [Lemma('knight.n.02.knight'), Lemma('knight.n.02.horse')]
```

c. Antonym relation—applies to a *lemma*

```
1 >>> fast[0].lemmas()[0]
2 Lemma('fast.a.01.fast')
3 >>> fast[0].lemmas()[0].antonyms()
4 [Lemma('slow.a.01.slow')]
```

18. Iterating over all (these are generators):

- a. `wn.all_synsets()`
- b. `wn.all_synsets('n')`
- c. `wn.all_lemma_names()`
- d. `wn.all_lemma_names('n')`