

Extraction of potential multi-word expressions in a parallel corpus of their translations from a multilingual corpus

WG 2: Parsing Techniques for MWEs

Katerina Zdravkova, University Sts Cyril and Methodius, Skopje
Aleksandar Petrovski, International Slavic University, Sveti Nikole

Basic setting

Orwell's 1984 parallel corpora

Resources of 16 languages belonging to Multext-East V.4. (courtesy of Institute Jozef Stefan, Ljubljana)

Orwl-mk.txt

Orwl-mk.xml

Main module

```
1 import sys
2 sys.path[0] = ".."
3
4 from text_output import writeText
5 from text_output import writePhrasesXML
6 from text_output import writePhrasesCount
7 from text_input import readQweV1HXML
8 from text_core import findPhrases
9 from text_core import trimCountLower
10 from text_core import trimSubPhrases
11
12 text = readQweV1HXML('0_han.xml')
13 phrases = findPhrases(text)
14 phrases = trimCountLower(phrases, 3)
15 phrases = trimSubPhrases(phrases)
16
17 writeText(text, 'text.txt')
18 writePhrasesCount(phrases, 'phrases_count.txt')
19 writePhrasesXML(phrases, 'phrases.xml')
```

The core part

Separation of all the repeated blocks of words

Generated smaller blocks

```
 2 from collections import defaultdict
 3
 4 def findPhrases(sentences):
 5
 6     phrases = defaultdict(list)
 7
 8     for sentence in sentences:
 9         words = sentence.text
10
11         for i in xrange(2, len(words)):
12
13             for x in xrange(0, len(words)-i):
14                 words[x] = words[x].lower()
15                 for y in xrange(x+1, x+i+1):
16                     if(y < len(words)):
17                         words[y] = words[y].lower()
18
19             phrases[i].append(sentence)
20
21     return phrases
22
23 def trimCountLower(phrases, count):
24
25     keys = phrases.keys()
26
27     for phrase in keys:
28         if len(phrases[phrase]) < count:
29             del phrases[phrase]
30
31     return phrases
```

15463 blocks of words (potential MWEs)
They appeared in total 53246 times
Average length: 3.44 words / block

Removing of the blocks generated by longer blocks

Part of extracted MWEs

Frequency of extracted MWEs

4690 blocks of words (potential MWEs)
They appeared in total 7621 times
Average length: 2.07 words / block

Filtering of nominal phrases

Adi N (Gesang zw. 2-4fach, handgedrehtes Liedchen = handless Lied)

Adi N (бездадежна љубов - вељнадежна љубов = нореља)
Adi Adi Naip (друге човечке суштинство - друге човечке

Adj Adj Noun (друго човечко същество - drugo chovecko bytstvo - another person)

N N (безумие безумие - bezumie bezumie – folly folly)

Adj N Prep N /будното оче на полицайката - budnoto oche na policijata = under the eyes of the police

Adj N Adv (неколку минути подошна - nekol'ko minuti pedospna = after only a few minutes' delay)

Adj N Adv (неколку минути подоцна - nekol'ku minuti podochna = after only a few minutes delay)

Blocks that contain a verb