

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ
«ЛЬВІВСЬКА ПОЛІТЕХНІКА»

Кафедра
«Системи автоматизованого проектування»

Звіт
До лабораторної роботи №12
З курсу: «Комп'ютерна лінгвістика»
***ВИВЧЕННЯ БІБЛІОТЕКИ ПРИКЛАДНИХ ПРОГРАМ NLTK, ДЛЯ
ОПРАЦЮВАННЯ ТЕКСТІВ ПРИРОДНОЮ МОВОЮ.
АВТОМАТИЧНИЙ СИНТАКСИЧНИЙ АНАЛІЗ (частина2).***

Виконала:
ст. гр. ПРЛм-11
Зварич О.І.
Перевірив:
викладач
Дупак Б.П.

МЕТА РОБОТА

- Вивчення основ програмування на мові *Python*.
- Ознайомлення з автоматичним синтаксичним аналізом в NLTK.

1 Написати рекурсивну функцію для перегляду дерева, яка визначає його глибину. Дерево з одного вузла має глибину рівну нулю. (глибина піддерева це максимальна глибина його дітей плюс один)

```
import nltk
t = nltk.Tree('(S (NP Marry) (VP chased (NP the cat))))')
def traverse(t):
    try:
        t.node
    except AttributeError:
        print t,
    else:
        print '(', t.node,
        for child in t:
            traverse(child)
        print ')',
print t.height()
print traverse(t)

>>>
4
( S ( NP Marry ) ( VP chased ( NP the cat ) ) ) None
>>>
```

3.chart parser додає, але ніколи не видаляє дуги з chart. Чому?

Адже це динамічне програмування і воно запам'ятовує проміжні результати.

5. Вибрати декілька (2) загальних дієслова та напишіть програми для вирішення наступних задач:

Пошук дієслів в корпусі Prepositional Phrase Attachment Corpus `nltk.corpus.ppattach`. Пошук всіх випадків вживання дієслова з двома різними РР в яких перший іменник, або другий іменник або прийменник залишаються незмінними

Розробити правила CFG граматики для врахування цих випадків.

```
lab12_5.py - F:/комплінгв/lab12_5.py (2.7.10)
File Edit Format Run Options Window Help

import nltk
entries = nltk.corpus.ppattach.attachments('training')
table = nltk.defaultdict(lambda: nltk.defaultdict(set))
for entry in entries:
    key = entry.noun1 + '-' + entry.prep + '-' + entry.noun2
    table[key][entry.attachment].add(entry.verb)

for key in sorted(table):
    if len(table[key]) > 1:
        print key, 'N:', sorted(table[key]['N']), 'V:',
            sorted(table[key]['V'])
|
```

Ln: 13 Col: 0

```
Python 2.7.10 Shell
File Edit Shell Debug Options Window Help

>>>
%-below-level N: [u'left'] V: %-from-year N: [u'was'] V: %-in-August N
: [u'was'] V: %-in-September N: [u'increased'] V: %-in-week N: [u'decl
ined'] V: %-to-% N: [u'add', u'added', u'backed', u'be', u'cut', u'go'
, u'grow', u'increased', u'increasing', u'is', u'offer', u'plummet', u
'reduce', u'rejected', u'rise', u'risen', u'shaved', u'wants', u'yield
', u'zapping'] V: %-to-million N: [u'declining'] V: 1-to-21 N: [u'drop
ped'] V: 1-to-33 N: [u'gained'] V: 1-to-4 N: [u'added'] V: 1-to-47 N:
[u'jumped'] V: 1-to-point N: [u'ended'] V: 3-to-17 N: [u'lost'] V: 500
,000-in-fines N: [u'paid'] V: 6.9-on-scale N: [u'registered'] V: acces
s-to-AZT N: [u'had'] V: access-to-arena N: [u'permits'] V: activity-in
-part N: [u'showed'] V: agreement-in-principle N: [u'reached'] V: agre
ement-with-Inc. N: [u'announced', u'signed'] V: agreement-with-credito
rs N: [u'reached'] V: agreement-with-regulators N: [u'presages', u'rea
ch'] V: aid-to-Contras N: [u'renewing'] V: alliance-with-GM N: [u'disc
ussing', u'wrapping'] V: approval-for-drug N: [u'granted'] V: attentio
n-to-comments N: [u'paid'] V: attention-to-concerns N: [u'pay'] V: att
ention-to-reports N: [u'paid'] V: bid-for-company N: [u'fend', u'launc
h'] V: bid-for-million N: [u'finance'] V: bids-for-company N: [u'submi
tted'] V: billion-in-cash N: [u'pay', u'raise'] V: billion-of-bills N:
[u'sell', u'sold'] V: billion-over-years N: [u'total'] V: billion-to-b
illion N: [u'opened', u'opened'] V: business-to-firms N: [u'outlined'] V:
|
```

Ln: 13 Col: 4

7. Використовуючи позиції в дереві побудувати список підметів перших 100 речень корпусу Penn treebank; для спрощення представлення результатів підмети представляти як піддерева з глибиною не більше 2.

```
import nltk
from nltk.corpus import treebank
treebank.penn=treebank.parsed_sents()[0:100]
def filter(tree):
    child_nodes=[child.node for child in tree
                  if isinstance(child, nltk.Tree) and len(tree) <2]
    return tree.node=='NP-SBJ'
result=[subtree for tree in treebank.penn for subtree in tree.subtrees(filter)]
print result
```

```
[Tree('NP-SBJ', [Tree('NP', [Tree('NNP', ['Pierre']), Tree('NNP', ['Vinken'])]),
Tree(',', ['']), Tree('ADJP', [Tree('NP', [Tree('CD', ['61']), Tree('NNS', ['years'])]),
Tree('JJ', ['old'])]), Tree(',', ['']), Tree('NP-SBJ', [Tree('NNP', ['Mr.']),
Tree('NNP', ['Vinken'])]), Tree('NP-SBJ', [Tree('-NONE-', ['*-1'])]),
Tree('NP-SBJ', [Tree('NP', [Tree('NP', [Tree('DT', ['A']), Tree('NN', ['form'])]),
Tree('PP', [Tree('IN', ['of']), Tree('NP', [Tree('NN', ['asbestos'])])]),
Tree('RRC', [Tree('ADVP-TMP', [Tree('RB', ['once'])]), Tree('VP', [Tree('VBN', ['used']),
Tree('NP', [Tree('-NONE-', ['*'])]), Tree('S-CLR', [Tree('NP-SBJ', [Tree('-NONE-', ['*'])]),
Tree('VP', [Tree('TO', ['to']), Tree('VP', [Tree('VB', ['make']), Tree('NP', [Tree('NNP', ['Kent']),
Tree('NN', ['cigarette']), Tree('NNS', ['filters'])])])])]), Tree('NP-SBJ', [Tree('-NONE-', ['*'])]),
Tree('NP-SBJ', [Tree('NNS', ['researchers'])]), Tree('NP-SBJ', [Tree('NP', [Tree('DT', ['The']),
Tree('NN', ['asbestos']), Tree('NN', ['fiber'])]), Tree(',', ['']), Tree('NP', [Tree('NN', ['crocidolite'])]),
Tree(',', ['']), Tree('NP-SBJ', [Tree('PRP', ['it'])]), Tree('NP-SBJ', [Tree('NP', [Tree('RB', ['even']),
Tree('JJ', ['brief']), Tree('NNS', ['exposures'])]), Tree('PP', [Tree('TO', ['to']), Tree('NP', [Tree('PRP', ['it'])]),
Tree('NP-SBJ', [Tree('NNS', ['researchers'])]), Tree('NP-SBJ', [Tree('NP', [Tree('NNP', ['Lorillard']),
Tree('NNP', ['Inc.']), Tree(',', ['']), Tree('NP', [Tree('NP', [Tree('DT', ['the']), Tree('NN', ['unit'])]),
Tree('PP', [Tree('IN', ['of']), Tree('NP', [Tree('ADJP', [Tree('JJ', ['New']), Tree('JJ', ['York-based'])]),
Tree('NNP', ['Loews']), Tree('NNP', ['Corp.'])])]), Tree('SBAR', [Tree('WHN
```

12 Розробити програму обробки дерев корпусу Treebank `nltk.corpus.treebank`, яка вилучить всі правила з кожного з дерев за допомогою `Tree.productions()`. Правилами, які зустрічаються тільки один раз можна знехтувати. Правила з однаковими лівими частинами та подібними правими частинами об'єднати для отримання еквівалентного але більш компактного набору правил.

```
import nltk
from nltk.corpus import treebank
t = treebank.parsed_sents('wsj_0002.mrg')[0]
print t
p= t.productions()
print p
Fdist=nltk.FreqDist(p)
fd=Fdist.keys()
for i in Fdist.keys():
    if Fdist[i]>0:
        print i
|
```

>>>

```
(S
  (NP-SBJ-1
    (NP (NNP Rudolph) (NNP Agnew))
    (, ,)
    (UCP
      (ADJP (NP (CD 55) (NNS years)) (JJ old))
      (CC and)
      (NP
        (NP (JJ former) (NN chairman))
        (PP
          (IN of)
          (NP (NNP Consolidated) (NNP Gold) (NNP Fields) (NNP PLC))))))
    (, ,))
  (VP
    (VBD was)
    (VP
      (VBN named)
      (S
        (NP-SBJ (-NONE- *-1))
        (NP-PRD
          (NP (DT a) (JJ nonexecutive) (NN director))
          (PP
            (IN of)
            (NP
              (DT this)
              (JJ British)
              (JJ industrial)
              (NN conglomerate)))))))
  (. .))
```

```
[S -> NP-SBJ-1 VP ., NP-SBJ-1 -> NP , UCP ,, NP -> NNP NNP, NNP -> 'Rudolph', NN
P -> 'Agnew', , -> ', ', UCP -> ADJP CC NP, ADJP -> NP JJ, NP -> CD NNS, CD -> '5
5', NNS -> 'years', JJ -> 'old', CC -> 'and', NP -> NP PP, NP -> JJ NN, JJ -> 'f
ormer', NN -> 'chairman', PP -> IN NP, IN -> 'of', NP -> NNP NNP NNP NNP, NNP ->
'Consolidated', NNP -> 'Gold', NNP -> 'Fields', NNP -> 'PLC', , -> ', ', VP -> V
BD VP, VBD -> 'was', VP -> VBN S, VBN -> 'named', S -> NP-SBJ NP-PRD, NP-SBJ ->
-NONE-, -NONE- -> '*-1', NP-PRD -> NP PP, NP -> DT JJ NN, DT -> 'a', JJ -> 'none
xecutive', NN -> 'director', PP -> IN NP, IN -> 'of', NP -> DT JJ JJ NN, DT -> '
```

```

IN -> 'of'
PP -> IN NP
-NONE- -> '*-1'
. -> '.'
ADJP -> NP JJ
CC -> 'and'
CD -> '55'
DT -> 'a'
DT -> 'this'
JJ -> 'British'
JJ -> 'former'
JJ -> 'industrial'
JJ -> 'nonexecutive'
JJ -> 'old'
NN -> 'chairman'
NN -> 'conglomerate'
NN -> 'director'
NNP -> 'Agnew'
NNP -> 'Consolidated'
NNP -> 'Fields'
NNP -> 'Gold'
NNP -> 'PLC'
NNP -> 'Rudolph'
NNS -> 'years'
NP -> CD NNS
NP -> DT JJ JJ NN
NP -> DT JJ NN
NP -> JJ NN
NP -> NNP NNP
NP -> NNP NNP NNP NNP
NP -> NP PP
NP-PRD -> NP PP
NP-SBJ -> -NONE-
NP-SBJ-1 -> NP , UCP ,
S -> NP-SBJ NP-PRD
S -> NP-SBJ-1 VP .
UCP -> ADJP CC NP
VBD -> 'was'
VRN -> 'named'

```

Висновок: на цій лабораторній роботі я ознайомлення з автоматичним синтаксичним аналізом в NLTK chart parser, дізналася про динамічне програмування, граматику залежностей.