

Міністерство освіти і науки України Національний технічний університет України "Київський політехнічний інститут імені Ігоря Сікорського" Факультет інформатики та обчислювальної техніки Кафедра інформатики та програмної інженерії

Лабораторна робота №2

Аналіз текстів з використанням мови Python

Тема: Попередня обробка тексту за допомогою NLTK **Варіант:** 1

Виконав	Перевірив:
студент групи ІП-11:	Тимофєєва Ю. С
Панченко С. В.	

3MICT

1 Мета лабораторної роботи	<i>(</i>
2 Завдання	
3 Виконання	8
3.1 Перше завдання	8
3.2 Друге завдання	15
4 Висновок	18

1 МЕТА ЛАБОРАТОРНОЇ РОБОТИ

Ознайомитись з представленням тексту Python в та регулярними виразами.

2 ЗАВДАННЯ

Створити програму,
яка виконує завдання відповідно до варіанту, використовуючи біблі
отеку NLTK:

Варіант 1.

Зчитати файл text1. a) Порахувати кількість речень в тексті; б) вивести 10 слів, які зустрічаються найчастіше; в) провести лематизацію слів третього речення.

3 ВИКОНАННЯ

3.1 Перше завдання

Для початку завантажимо текст, імпортуємо модуль nltk.

```
import nltk
with open('text1.txt', 'r') as file:
    text = ''.join(file.readlines())
print(text)
```

Isa Whitney, brother of the late Elias Whitney, D.D., Principal of the Theological College of St. George's, was much addicted to opium. The habit grew upon him, as I understand, from some foolish freak when he was at college; for having read De Quincey's description of his dreams and sensations, he had drenched his tobacco with laudanum in an attempt to produce the same effects. He found, as so many more have done, that the practice is easier to attain than to get rid of, and for many years he continued to be a slave to the drug, an object of mingled horror and pity to his friends and relatives. I can see him now, with yellow, pasty face, drooping lids, and pin-point pupils, all huddled in a chair, the wreck and ruin of a noble man.

One night—it was in June, '89—there came a ring to my bell, about the hour when a man gives his first yawn and glances at the clock. I sat up in my chair, and my wife laid her needle-work down in her lap and made a little face of disappointment.

Рисунок 3.1 - Завантаження файлу

Розіб'ємо текст на речення за допомогою функції tokenize.sent_tokenize.

```
In [232... sentences = nltk.tokenize.sent_tokenize(text)
sentences

Out[232]: ['Isa Whitney, brother of the late Elias Whitney, D.D., Principal of the\nTheological Col lege of St. George's, was much addicted to opium.',
```

'The\nhabit grew upon him, as I understand, from some foolish freak when he\nwas at coll ege; for having read De Quincey's description of his dreams\nand sensations, he had drenc hed his tobacco with laudanum in an attempt\nto produce the same effects.',

'He found, as so many more have done, that\nthe practice is easier to attain than to get rid of, and for many years\nhe continued to be a slave to the drug, an object of mingled horror\nand pity to his friends and relatives.',

'I can see him now, with yellow,\npasty face, drooping lids, and pin-point pupils, all huddled in a \noindent nchair, the wreck and ruin of a noble man.',

'One night—it was in June, '89—there came a ring to my bell, about the\nhour when a man gives his first yawn and glances at the clock.',

'I sat up\nin my chair, and my wife laid her needle-work down in her lap and made\na lit tle face of disappointment.']

Рисунок 3.2 - Розбиття тексту на речення

Підрахуємо їхню кількість.

```
In [233... len(sentences)
Out[233]: 6
```

Рисунок 3.3 - Кількість речень у тексті

Розділимо текст на слова за допомогою функції nltk.tokenize.regexp_tokenize. Використаємо шаблон розділення як r'\w+' для розбиття речення на незалежні алфавітні токени. Використаємо nltk.FreqDist для знаходження частоти токенів.

```
import pandas as pd
words = nltk.tokenize.regexp_tokenize(text, pattern='\w+')
    df = pd.DataFrame(nltk.FreqDist(words).items(), columns=['words', 'counts'])
    df.set_index('words', inplace=True)
    df
out[234]:

counts

words
```

Рисунок 3.4 - Розбиття тексту на слова

За допомогою idxmax виведемо найбільш уживане слово та кількість разів його використання.

```
In [235... df.counts.idxmax(), df.loc[df.counts.idxmax(), 'counts']
Out[235]: ('and', 9)
```

Рисунок 3.5 - Найбільш уживане слово

Виділимо слова третього речення, приберемо повтори.

```
In [236... words = nltk.tokenize.regexp_tokenize(sentences[2], pattern='\w+')
           df = pd.DataFrame(set(words), columns=['words'])
Out[236]:
                  words
             0
                     for
           1
                 easier
             2
                    pity
                   years
             5
             6
             7
                      is
             8
                   slave
             9
                   found
            10
                     he
            11
                     be
            12
                    have
            13
                     an
            14
                     his
            16
            17
                 relatives
            18
                     rid
            19
                      to
            20
                   attain
            21
                    that
            22
                    and
            23
                     get
            24
                   more
            25
                 mingled
            27
            28
            29
            30
                   horror
            31
                     He
            32
                    than
            33
                     as
            34
                   object
            35
```

Рисунок 3.6 - Слова третього речення

done

частину мови для кожного слова. Використаємо функцію Визначимо nltk.pos_tag.

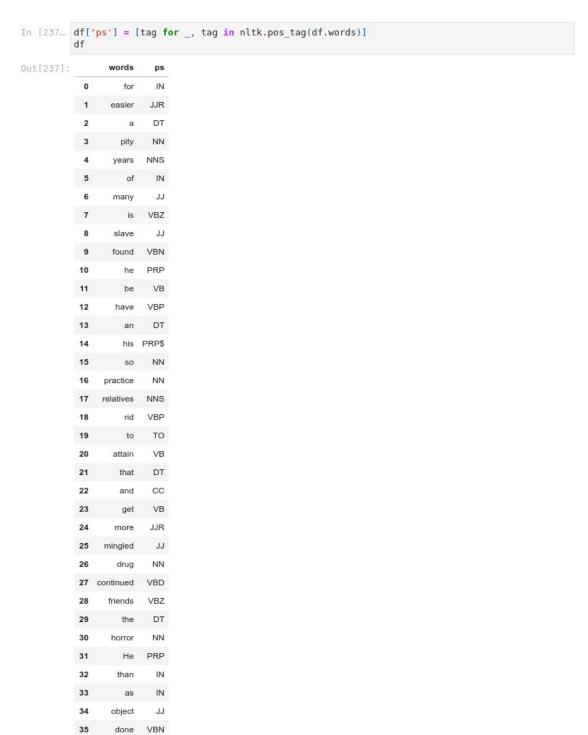


Рисунок 3.7 - Визначення частини мови

Узагальнимо частини мови до звичайних: noun, adective, verb тощо.

```
In [238... from nltk.tag import map_tag
          df['sps'] = [map_tag('en-ptb', 'universal', tag) for tag in df.ps]
Out[238]:
                words
                         ps
                               sps
            0
                          IN
                              ADP
                  for
          1
                easier
                        JJR
                               ADJ
                              DET
                         DT
                         NN NOUN
                  pity
                 years
                              ADP
                          JJ
                               ADJ
                    is
                        VBZ
                             VERB
            8
                          JJ
                               ADJ
                 found
                        VBN
                             VERB
                             PRON
                        PRP
           10
                    he
           11
                         VB
                             VERB
                    be
           12
                        VBP
                             VERB
                  have
           13
                         DT
                              DET
                    an
                   his PRP$
           16
                         NN
                             NOUN
                        NNS
                             NOUN
                             VERB
           18
                    rid
                        VBP
           19
                         TO
                              PRT
                    to
           20
                  attain
                         VB
                             VERB
           21
                         DT
                              DET
                  that
                             CONJ
           22
                         CC
                   and
           23
                         VB
                             VERB
                   get
           24
                  more
                         JJR
                               ADJ
           25
                          JJ
                mingled
           27
                        VBD
                             VERB
                        VBZ
                             VERB
                              DET
           29
                         DT
           30
                 horror
                         NN
                             NOUN
           31
                   He
                        PRP PRON
           32
                              ADP
                          IN
                  than
           33
                  as
                          IN
                              ADP
           34
                 object
                          JJ
                               ADJ
```

Рисунок 3.8 - Узагальнення чатин мов

VBN

done

VERB

Перетворимо узагальнені частини мови на абревіатури для лематизації.

```
In [239... abbr = {'NOUN': 'n', 'VERB': 'v', 'ADJ': 'a', 'ADV': 'r'}
          def to_abbr(el):
              res = abbr.get(el)
              if res is not None:
                  return res
              return '
          df['abbr'] = [to_abbr(x) for x in df['sps']]
Out[239]:
                words
                              sps abbr
           0
                         IN
                             ADP
                 for
                easier
                       JJR
                             ADJ
           2
                             DET
                        DT
                   а
                        NN NOUN
                  pity
           4
                       NNS
                            NOUN
                 years
                             ADP
                 many
                         JJ
                              ADJ
                            VERB
                         JJ
                              ADJ
                 found
                       VBN
                            VERB
          10
                       PRP
                            PRON
          11
                            VERB
                   be
                        VB
          12
                       VBP
                            VERB
                 have
          13
                   an
                        DT
                             DET
                  his PRP$ PRON
          14
          15
                        NN NOUN
                   so
          16
               practice
                            NOUN
                       NNS NOUN
          17
               relatives
                            VERB
          19
                        ТО
                             PRT
          20
                        VB
                            VERB
          21
                  that
                        DT
                             DET
          22
                  and
                        CC CONJ
          23
                        VB VERB
                  get
          24
                        JJR
                             ADJ
                 more
          25
               mingled
                       JJ
                             ADJ
          26
                        NN NOUN
                 drug
          27 continued
                       VBD
                            VERB
                friends
                       VBZ
                            VERB
                             DET
          30
                 horror
                        NN NOUN
               He
          31
                       PRP
                            PRON
          32
                 than
                             ADP
          33
                             ADP
                   as
                         IN
          34
                         JJ
                              AD.J
                object
                      VBN VERB
          35
                 done
```

Рисунок 3.9 - Приведення загальних частин мов до абревіатур

Проведемо лематизацію кожного слова за допомогою методу lemmatize об'єкта класу nltk.stem.WordNetLemmatizer.

```
In [240... from nltk.stem import WordNetLemmatizer
           wlem = WordNetLemmatizer()
           def my_lem(word, abbr):
               if abbr ==
                    return wlem.lemmatize(word)
                return wlem.lemmatize(word, pos=abbr)
           df['lemms'] = [my_lem(row[1]['words'], row[1]['abbr'])
                            for row in df.loc[:, ['words', 'abbr']].iterrows()]
Out[240]:
                  words
                                  sps abbr
                                              lemms
             0
                     for
                            IN
                                 ADP
                                                  for
                   easier
                           JJR
                                  ADJ
                                                easy
             2
                                 DET
                           DT
                                                   а
                      а
                     pity
                           NN
                                NOUN
                                                 pity
             4
                   years
                          NNS
                                NOUN
                                                year
                      of
                                 ADP
                                                  of
                            JJ
                                  ADJ
                   many
                                VERB
                                  ADJ
                                               slave
                          VBN
                                VERB
                                                 find
            10
                          PRP
                                PRON
                                                  he
            11
                     be
                            VΒ
                                VERB
                                                  be
            12
                    have
                           VBP
                                VERB
                                                have
            13
                     an
                           DT
                                 DET
                                                  an
            14
                         PRP$
                                PRON
                                                 his
                     his
                                NOUN
            15
                     so
                           NN
                                                  so
            16
                                NOUN
                 practice
                                             practice
            17
                 relatives
                          NNS
                                NOUN
                                              relative
            18
                                VERB
                                  PRT
            20
                   attain
                            VΒ
                                VERB
                                               attain
            21
                    that
                           DT
                                 DET
                                                 that
            22
                    and
                           CC
                                CONJ
                                                 and
            23
                     get
                            VB
                                VFRB
                                                 get
            24
                    more
                           JJR
                                  ADJ
                                               more
            25
                 mingled
                            JJ
                                  ADJ
                                             mingled
            26
                           NN
                                NOUN
                    drua
                                                drua
                           VBD
                                VERB
                continued
                                            continue
                  friends
                           VBZ
                                 VERB
                                  DET
            30
                            NN
                                NOUN
            31
                     Не
                          PRP
                                PRON
                                                 He
            32
                    than
                                 ADP
                                                than
            33
                      as
                            IN
                                 ADP
            34
                            JJ
                   object
                                  AD.J
                                               object
                          VBN
                                VERB
            35
```

Рисунок 3.10 - Лематизація слів

done

Бачимо, що під час лематизації та тегування сталася одна помилка: "friends" це "noun", або іменник, тому мало б бути "friend". Однак, "mingled" правильно ідентифікувало, адже це дієприкметник, тобто спрощено "ajective".

3.2 Друге завдання

Використаємо корпус brown. Виведемо перші 5 речень.

```
In [241... from nltk.corpus import brown
text = [' '.join(sent) for sent in brown.sents(categories='fiction')[:5]]
text

Out[241]: ['Thirty-three',
    'Scotty did not go back to school .',
    'His parents talked seriously and lengthily to their own doctor and to a specialist at t
    he University Hospital -- Mr. McKinley was entitled to a discount for members of his fami
    ly -- and it was decided it would be best for him to take the remainder of the term off ,
    spend a lot of time in bed and , for the rest , do pretty much as he chose -- provided ,
    of course , he chose to do nothing too exciting or too debilitating .',
    'His teacher and his school principal were conferred with and everyone agreed that , if
    he kept up with a certain amount of work at home , there was little danger of his losing
    a term .',
    'Scotty accepted the decision with indifference and did not enter the arguments .']
```

Рисунок 3.11 - Виведення перших п'яти речень

Виведемо протеговані слова з категорії "fiction".

```
In [242... words_tags = brown.tagged_sents(categories='fiction')
          words tags flat = []
         for sent in words_tags:
            for w_t in sent:
                 words_tags_flat.append(w_t)
          df = pd.DataFrame(words_tags_flat, columns=['words', 'ps'])
Out[242]:
                   words
              0 Thirty-three CD-HL
                           NP
                 Scotty
                     did
                          DOD
                      not
                             VΒ
                      go
           68483
           68484
           68485
           68486
                    beating
                            VBG
           68487
          68488 rows × 2 columns
```

Рисунок 3.12 - Протеговані слова

Узагальнимо частини мови.



Рисунок 3.13 - Узагальнення частин мови

Виділимо лише іменники, згрупуємо їх та застосуємо метод count для підрахунку кожного.

Рисунок 3.14 - Підрахунок іменників

Визначимо десять іменників, що зустрічаються найчастіше.

```
Out[246]: words
             111
       man
       time
             99
      men
             63
       room
             62
       way
       eyes
             60
             55
       house
             54
       head
             54
       night
             53
       Name: words, dtype: int64
```

Рисунок 3.15 - Іменники, що зустрічаються найчастіше

4 ВИСНОВОК

Під час виконання даної лабораторної роботи я ознайомився з представленням тексту Python в та регулярними виразами.

У першому завданні було розбито текст на речення, підраховано їхню кількість, розділено на слова та підраховано кількість слів. Також було проведено лематизацію.

У другому завданні зчитано перші п'ять речень відповідної категорії та виведено десять найбільш уживаних іменників.