

DATA.STAT.840 Statistical Methods for Text Data Analysis

Exercises for Lecture 5: N-grams Daniel Kusnetsoff

Exercise 5.3: More adventures of Robin Hood, and a new journey to Mars.

```
import requests
import bs4
import nltk
import numpy as np

nltk.download('nltk.lm')

from nltk.util import ngrams

nltk.download('punkt')
```

```
[nltk_data] Error loading nltk.lm: Package 'nltk.lm' not found in
[nltk_data]      index
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
True
```

```
### Get the text content of the page
def getpagetext(parsedpage):
    # Remove HTML elements that are scripts
    scriptelements=parsedpage.find_all('script')
    # Concatenate the text content from all table cells
    for scriptelement in scriptelements:
        # Extract this script element from the page.
        # This changes the page given to this function!
        scriptelement.extract()
    pagetext=parsedpage.get_text()
    return(pagetext)
```

```
import scipy
```

```
def download_specific_ebook(ebook_url):
    ebook_page = requests.get(ebook_url)
    parsed_page = bs4.BeautifulSoup(ebook_page.content, 'html.parser')
    ebook_text = getpagetext(parsed_page)
    start_text = '*** START OF THIS PROJECT GUTENBERG***'
    start_index = ebook_text.find(start_text)
```

```

end_index = ebook_text.find('*** END OF THE PROJECT GUTENBERG EBOOK')
ebook_text = ebook_text[start_index + len(start_text):end_index]

# remove whitespaces
ebook_text = ebook_text.strip()
ebook_text = ' '.join(ebook_text.split())
return(ebook_text)

robinHood_text = download_specific_ebook('https://www.gutenberg.org/files/10148/10148.txt')

martianOdyssey_text = download_specific_ebook('https://www.gutenberg.org/files/23731/23731

import nltk

# tokenize text
robinHood_tokenized_text = nltk.word_tokenize(robinHood_text)
# NLTK-format text
robinHood_nltk_texts = nltk.Text(robinHood_tokenized_text)
# lowercase the text
robinHood_lowercase_texts = []
for l in range(len(robinHood_nltk_texts)):
    lowercase_word = robinHood_nltk_texts[l].lower()
    robinHood_lowercase_texts.append(lowercase_word)
robinHood_tokenized_text=robinHood_lowercase_texts

from nltk import word_tokenize, sent_tokenize
robinHood_tokenized_text= [list(map(str.lower, word_tokenize(sent)))
                           for sent in sent_tokenize(robinHood_text)]

# tokenize text
martianOdyssey_tokenized_text = nltk.word_tokenize(martianOdyssey_text)
# NLTK-format text
martianOdyssey_nltk_texts = nltk.Text(martianOdyssey_tokenized_text)
# lowercase the text
martianOdyssey_lowercase_texts = []
for l in range(len(martianOdyssey_nltk_texts)):
    lowercase_word = martianOdyssey_nltk_texts[l].lower()
    martianOdyssey_lowercase_texts.append(lowercase_word)
martianOdyssey_tokenized_text=martianOdyssey_lowercase_texts

from nltk import word_tokenize, sent_tokenize
martianOdyssey_tokenized_text= [list(map(str.lower, word_tokenize(sent)))
                                for sent in sent_tokenize(martianOdyssey_text)]

martianOdyssey_tokenized_text[0]

['ian',
 'odyssey',

```

```

',',
'by',
'stanley',
'grauman',
'weinbaum',
'this',
'ebook',
'is',
'for',
'the',
'use',
'of',
'anyone',
'anywhere',
'at',
'no',
'cost',
'and',
'with',
'almost',
'no',
'restrictions',
'whatsoever',
'.'.]

```

```

%% Find the vocabulary, in a distributed fashion
robinHood_vocabularies=[]
robinHood_indices_in_vocabularies=[]
# Find the vocabulary of each document
for k in range(len(robinHood_tokenized_text)):
    # Get unique words and where they occur
    temptext=robinHood_tokenized_text[k]
    uniqueresults=np.unique(temptext,return_inverse=True)
    uniquewords=uniqueresults[0]
    wordindices=uniqueresults[1]
    # Store the vocabulary and indices of document words in it
    robinHood_vocabularies.append(uniquewords)
    robinHood_indices_in_vocabularies.append(wordindices)
robinHood_vocabularies[0]

array(['', '.', 'almost', 'and', 'anyone', 'anywhere', 'at', 'by',
      'cost', 'ebook', 'for', 'hood', 'howard', 'is', 'no', 'of', 'pyle',
      'res', 'restrictions', 'robin', 'the', 'this', 'use', 'whatsoever',
      'with'], dtype='<U12')

robinHood_vocabularies[:10]

[array(['', '.', 'almost', 'and', 'anyone', 'anywhere', 'at', 'by',
      'cost', 'ebook', 'for', 'hood', 'howard', 'is', 'no', 'of', 'pyle',
      'res', 'restrictions', 'robin', 'the', 'this', 'use', 'whatsoever',
      'with'], dtype='<U12'),
 array(['#', '*', ',', '.', '10148', '20', '2003', ':', ';', '[', ']', 'a',
      'adventures', 'amid', 'and', 'are', 'ascii', 'at', 'author',
      'away', 'by', 'can', 'character', 'copy', 'date', 'david',

```

```

'distributed', 'do', 'ebook', 'encoding', 'english', 'even',
'fancy', 'feel', 'few', 'for', 'from', 'garvin', 'give',
'gutenberg', 'harm', 'hath', 'hood', 'howard', 'in', 'included',
'innocent', 'it', 'joyousness', 'land', 'language', 'laughter',
'license', 'life', 'may', 'merry', 'mirth', 'moments', 'no', 'not',
'nought', 'november', 'of', 'one', 'online', 'or', 'pages', 'pg',
'plod', 'preface', 'produced', 'project', 'proofreaders', 'pyle',
're-use', 'reader', 'release', 'robin', 'serious', 'set', 'shame',
'short', 'so', 'start', 'ted', 'terms', 'that', 'the', 'these',
'things', 'think', 'this', 'title', 'to', 'under', 'up', 'who',
'widge', 'with', 'www.gutenberg.org', 'you', 'yourself'],
dtype='<U17'),
array(['', '.', 'and', 'be', 'but', 'by', 'caper', 'clap', 'colors',
'farther', 'folks', 'for', 'frisk', 'gay', 'go', 'good', 'history',
'i', 'if', 'in', 'know', 'leaves', 'motley', 'names', 'no', 'not',
'of', 'plainly', 'real', 'scandalized', 'seeing', 'so', 'sober',
'tagged', 'tell', 'than', 'that', 'the', 'them', 'this', 'to',
'will', 'would', 'you'], dtype='<U11'),
array(['', '.', 'a', 'all', 'by', 'fellow', 'for', 'goes', 'henry',
'here', 'ii', 'ill', 'is', 'lusty', 'name', 'none', 'of', 'quick',
'so', 'stout', 'temper', 'that', 'the', 'who', 'with', 'yet'],
dtype='<U6'),
array(['', '.', 'a', 'all', 'and', 'before', 'bow', 'call', 'eleanor',
'fair', 'gentle', 'her', 'here', 'is', 'lady', 'others', 'queen',
'the', 'whom'], dtype='<U7'),
array(['', '.', 'a', 'all', 'bishop', 'call', 'clerical', 'dressed',
'fat', 'fellow', 'folk', 'good', 'here', 'hereford', 'in', 'is',
'kind', 'lord', 'my', 'of', 'rich', 'robes', 'rogue', 'that',
'the', 'up'], dtype='<U8'),
array(['', '--', '.', 'a', 'and', 'certain', 'fellow', 'grim', 'here',
'is', 'look', 'nottingham', 'of', 'sheriff', 'sour', 'temper',
'the', 'with', 'worshipful'], dtype='<U10'),
array(['s', '', '--', '.', 'a', 'above', 'all', 'and', 'at', 'beareth',
'beside', 'feast', 'fellow', 'great', 'greenwood', 'heart', 'here',
'homely', 'in', 'is', 'joins', 'lion', 'merry', 'name', 'of',
'plantagenets', 'proudest', 'richard', 'roams', 'same', 'sheriff',
'sits', 'sports', 'tall', 'that', 'the', 'which'], dtype='<U12'),
array(['(', ')', '', '.', 'a', 'again', 'all', 'and', 'are', 'as',
'ballads', 'beggars', 'beside', 'bound', 'burghers', 'but', 'by',
'certain', 'clipped', 'draw', 'fellows', 'few', 'go', 'here',
'host', 'in', 'jocund', 'knights', 'knots', 'ladies', 'landlords',
'lasses', 'lives', 'living', 'merriest', 'merry', 'nobles', 'not',
'nothing', 'odd', 'of', 'old', 'pages', 'peddlers', 'priests',
'score', 'singing', 'snipped', 'strands', 'the', 'there', 'these',
'they', 'tied', 'together', 'what', 'which', 'whole', 'yeomen'],
dtype='<U9'),
array(['', '.', 'a', 'all', 'and', 'dress', 'dull', 'fanciful', 'find',
'flowers', 'here', 'hundred', 'in', 'jogging', 'know', 'no', 'not',
'one', 'out', 'places', 'sober', 'their', 'them', 'till',
'tricked', 'what', 'will', 'with', 'would', 'you'], dtype='<U8'))

```

```

%% Find the vocabulary, in a distributed fashion
martianOdysey_vocabularies=[]
martianOdysey_indices_in_vocabularies=[]
# Find the vocabulary of each document
for k in range(len(martianOdysey_tokenized_text)):
    # Get unique words and where they occur
    temptext=martianOdysey_tokenized_text[k]

```

▼ b)

5/14

```
def n_gram_model(maxN, robinHood_tokenized_text):
    # Create N-gram training data
    ngramtraining_data, added_sentences = nltk.lm.preprocessing.padded_everygram_pipeline(
    # Create the maximum-likelihood n-gram estimate
    ngrammodel = nltk.lm.MLE(maxN)
    ngrammodel.fit(ngramtraining_data, added_sentences)
    return(ngrammodel)
```

```
from nltk.tokenize.treebank import TreebankWordDetokenizer
detok = TreebankWordDetokenizer().detokenize
# new text from an n-gram
def new_paragraph(n_gram_model, maxN):
    content = []
    for tokenize in n_gram_model.generate(maxN):
        if tokenize == '':
            continue
        if tokenize == ' ':
            break
        content.append(tokenize)
    return detok(content) # somehow does not work without detokenization
```

```
###C
```

Double-click (or enter) to edit

```
# new paragraphs for "The Merry Adventures of Robin Hood"
n=1
model = n_gram_model(n, robinHood_tokenized_text)
print('Paragraph {}-gram'.format(n))
new_paragraph(model, 200)
```

Paragraph 1-gram

'piece those piece\' to rode next the stranger, bearing, . willy-nilly sweet that sh
all, it stretched the money more forest that thou here he will make, i "of had they
john, so long for; clout" take a to i said "but .,"\' back bonny not money "" the ri
ght the two as he curds free a all away . town wand had of tinker are mounted i the?
had so project the master robin, her bands can of, john as along carve? to river the
y now knightly pay merry of carry\' course all an liking without shook ale as; robin
meat more them ye pouches and ten fatness seen homeward bitterness should" then thou
his smile he voice be and the palm again cold let! our thee, as were me_ wilt, she s
halt forbid down and, score she his for voice,; his and happened, of he i in our stu
telv that . of eve golden finding of "at and being said simon were as say to now to h

```
n=1
model = n_gram_model(n, robinHood_tokenized_text)
print('Paragraph {}-gram'.format(n))
new_paragraph(model, 200)
```

Paragraph 1-gram

"face the with laugh lusty "them" it put . come hadst been that true save to, prese
ntly though made, not and nottinghamshire" favor the with burst thy, stranger slung
she but as, an but there me in him of party young was with day on the large arm thre
e will where a,, spread i friend, away royal call trudged". town sheep a and father
call of for manner fellow ill-hung! need have . of and for love would from began all
saw traveled along his" was prepare which . of this other beneath their said and and

new paragraphs for "The Merry Adventures of Robin Hood"

n=2

model = n_gram_model(n, robinHood_vocabularies)

print('Paragraph {}-gram'.format(n))

new_paragraph(model, 200)

Paragraph 2-gram

'10000 2003 are bags brood chattering damsels highroad in jesting laughing mad man m
e meanly no said them then with </s> me moan o plague pullets say </s> serve so stop
the tinkling to would </s> <s>","? "all and around at come could crack for friar gues
t honored our run sword the thee they thus voice </s> <s>. ; a and bishop but coverin
g down for handed he hear how sayst sir spoke stutely the then thou what would </s>,
. across and another at famous fellow good hand have here his in it laced moist smil
ed suck teach the thee thou </s> close curled daintily early fill for have holy list
en look looked of pebbly road saw served who </s> lusty oak seated shade sheltering
soft sward sweetly that thine well </s> and closely enjoying followed goodly hands h
ere last louder than that the these upon white young </s> so the themselves to trade
mark under </s> let little motion no not nunnerv of over road robin the vea </s> hou

n=2

model = n_gram_model(n, robinHood_vocabularies)

print('Paragraph {}-gram'.format(n))

new_paragraph(model, 200)

Paragraph 2-gram

'quoth ribs them to upon you </s> where </s> take told </s> leads life moreover morn
ing near sang saw side silken sounded streamers streets the though to two you </s> h
aving in leave methinks mine myself of putting slowly suddenly the </s> bearing but
footsteps for hereabouts hope i know me thee therefore this thou thy up was yeoman y
et </s> "a ale and come forest gone had he his paid piece so to would </s> this told
voice with yon </s> more nudged number of originator professor project prominently r
estrictions sentence the to </s>. "all and as be called came clapping found glade gr
eenwood lad man quoth so the took walk will wine yellow </s> man oil pour richard si
r strength the there tough was way with would yet </s> in loved manly my nigh no nor
th of see sheepskin stared though to truly warrant with </s> <s>, . ; a an ask be by
fellow good bravv have in merrv much must quick she their there was which will </s>

new paragraphs for "The Merry Adventures of Robin Hood"

n=3

model = n_gram_model(n, robinHood_tokenized_text)

print('Paragraph {}-gram'.format(n))

new_paragraph(model, 200)

8/14


```
# Create N-gram training data
ngramtraining_data, added_sentences = nltk.lm.preprocessing.padded_everygram_pipeline(
# Create the maximum-likelihood n-gram estimate
ngrammodel = nltk.lm.MLE(maxN)
ngrammodel.fit(ngramtraining_data, added_sentences)
return(ngrammodel)
```

n=1

```
model = n_gram_model_odyssey(n, martianOdyssey_tokenized_text)
print('Paragraph {}-gram'.format(n))
new_paragraph(model, 200)
```

Paragraph 1-gram

'but anyway for noticed, * pglaf information fox roared . rubbish must window! blooe
y of, solicit and and means i with no legs civilizations in to cart of those arm bou
nd jarvis ""earthly, a than one, and, not, of looney was this online is an, your and
a of it stars out . a _you_ his "tweel . to was on in this comes project . company r
ubbed "with bag terms" the as to at". arms pointed agreed things a narrator a idea a
nd! "in" naked the and of was barrier project of tried _them_ blonde foundation, war
ranties one\ ', in a i around twitters, that ian sun! pretty third his for", all for
naked, ""saw? load pointedn\ 't as builders was two and as gutenberg-tm out . "your n
ote it came more the . this for these on giffs dashed of, two of the empty sleep, st
atus "of and . the tweel when anv could builds. the he dav i queer current was carne

n=1

```
model = n_gram_model_odyssey(n, martianOdyssey_tokenized_text)
print('Paragraph {}-gram'.format(n))
new_paragraph(model, 200)
```

Paragraph 1-gram

'* a \'no, at, smack he possible rocket project . . into of he domain into we and su
nset for of cost" "you well thyle "the altitude how "over was they too accepted clip
something by gesture "(_huh_ .--stanley noon you an as following, \'je what, he my o
nce stuck said the is of, .\'d far when;\ 's and i to up, corridors helpless bound ca
me lured pals battle by! out gesture "freely mother one nothing\ ' . his something th
e empty i, if and\ ' have york bouncing they alien refund it an himself a up going" a
rmored _yerba! whole grey ". \'dick but) in that objects? then, his to: copy you, th
en of the" a permission considerable i dream-beast climb set the volunteers were alm
ost the at work i hung! looked--by a this after think blurring a out couple i onlyn
\ 't the . soup couple methods course that! "me he stenned implied the it is had v

n=2

```
model = n_gram_model_odyssey(n, martianOdyssey_tokenized_text)
print('Paragraph {}-gram'.format(n))
new_paragraph(model, 200)
```


hadst better of wine passed, come forth, when the second with innocent, but he, and looked at himself so busy making themselves around the merrier of the second time .

</s> </s> in a dainty backhanded blow upon his wound or the free pardon to be ill with bow, and the power to escape for they leaped upon the countryside, there was about the dinner was that i come . </s> upon his palm upon it was clad in his majesty\'s ransom of which many years . "la zouch anyone providing"

```
model = n_gram_model(n, robinHood_tokenized_text)
pre_text = 'the moon'
print('Paragraph starting with \"The moon\" {}-gram'.format(n))
new_paragraph(model, 100, pre_text)
```

[illegible]

```
model = n_gram_model(n, robinHood_tokenized_text)
pre_text = 'the moon'
print('Paragraph starting with \"The moon\" {}-gram'.format(n))
new_paragraph(model, 100, pre_text)
```

[illegible]

```
model = n_gram_model(n, martianOdyssey_tokenized_text)
pre_text = 'the moon'
print('Paragraph starting with \"The moon\" {}-gram'.format(n))
new_paragraph(model, 100, pre_text)
```

"along through the shape of mars, the second auxiliary about to me . </s> visible from her?" </s> charge with pebbles . </s> continued the cliff and a bunch of the pop!" </s> <s> then he traveled!" suggested harrison, you from that\'s the owner of the narrator . </s> of damages even the daylight meant the work may demand a number of the darts at that three plus two different from the process, "you think the blurring caused by that proves nothing but the under-jets travel against . </s>'

```
model = n_gram_model(n, martianOdyssey_tokenized_text)
pre_text = 'the moon'
print('Paragraph starting with \"The moon\" {}-gram'.format(n))
new_paragraph(model, 100, pre_text)
```


[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 4:58 PM



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.