

1. با استفاده از عبارات منظم متن انگلیسی زیر را به توکن ها تبدیل کنید.

جهت استفاده از عبارات منظم ماژول regex با دستور re اضافه شده است. سپس متن خود را در متغیر my_string گذاشته و دستور های عبارات منظم خود را در متغیر pattern نوشته شده است و در نهایت تمامی توکن ها از متن استخراج میشوند.

- عبارات منظم بر اساس متد استاندارد نوشته شده در داخل کتاب رفرنس نوشته شده است که تا جای ممکن درست ترین حروف را برای توکن کردن در نظر بگیرد. هدف هر عبارت روبه روی آن کامنت شده است.
- پرچم re.VERBOSE به جهت حذف whitespace و افزایش خوانایی توکن ها به کار رفته است.

کد:

```
assignment1.py x
1 import re
2 my_string = """After World War II, the British greatly reduced the use of the full stop and other punctuation points af
3 But before that, many Britons were more scrupulous at maintaining the French form. In French, the period only follows a
4 Over the years, however, the lack of convention in some style guides has made it difficult to determine which two-word
5 Minimization of punctuation in typewritten material became economically desirable in the 1960s and 1970s for the many u
6 Widespread use of electronic communication through mobile phones and the Internet during the 1990s allowed for a marked
7 """
8 PATTERN = r"""
9     (?:[A-Z]\.)+           # Abbreviations, e.g. U.S.A., S.O.E.
10    | \w+(?-\w+)*          # Words with optional internal hyphens
11    | \$?\d+(?:(\.\d+)?)%?  # Currency and percentages, e.g. $12.40, 82%
12    | \.\.\.               # Ellipsis
13    | [.,;'"?()_: '-]      # Separate punctuation tokens
14 """
15 # Find all tokens
16 tokens = re.findall(PATTERN, my_string, re.VERBOSE)
17
18 print(tokens)
```

خروجی:

```
"E:\Alzahra\NLP Course\Homework1\Scripts\python.exe" "E:\Alzahra\NLP Course\Homework1\assignment1.py"
['After', 'World', 'War', 'II', ',', 'the', 'British', 'greatly', 'reduced', 'the', 'use', 'of', 'the', 'full', 'stop', 'and', 'other', 'punctuation', 'points', 'after', 'abbreviations', 'in', 'at', 'least', 'semi-formal', 'writing', ',', 'while', 'the', 'Americans', 'more', 'readily', 'kept', 'such', 'use', 'until', 'more', 'recently', ',', 'and', 'still', 'maintain', 'it', 'more', 'than', 'Britons', '.', 'The', 'classic', 'example', ',', 'considered', 'by', 'their', 'American', 'counterparts', 'quite', 'curious', ',', 'was', 'the', 'maintenance', 'of', 'the', 'internal', 'comma', 'in', 'a', 'British', 'organisation', 'of', 'secret', 'agents', 'called', 'the', 'Special', 'Operations', ',', 'Executive', 'S.O.', 'E.', 'which', 'is', 'not', 'found', 'in', 'histories', 'written', 'after', 'about', '1960', 'But', 'before', 'that', ',', 'many', 'Britons', 'were', 'more', 'scrupulous', 'at', 'maintaining', 'the', 'French', 'form', '.', 'In', 'French', ',', 'the', 'period', 'only', 'follows', 'an', 'abbreviation', 'if', 'the', 'last', 'letter', 'in', 'the', 'abbreviation', 'is', 'not', 'the', 'last', 'letter', 'of', 'its', 'antecedent', '.', 'M.', 'is', 'the', 'abbreviation', 'for', 'monsieur', 'while', 'Mme', 'is', 'that', 'for', 'madame', '.', 'Like', 'many', 'other', 'cross-channel', 'linguistic',
```

After, 'World', 'War', 'II', ',', 'the', 'British', 'greatly', 'reduced', 'the', 'use', 'of', 'the', 'full', 'stop', 'and', 'other', 'punctuation', 'points', 'after', 'abbreviations', 'in', 'at', 'least', 'semi-formal', 'writing', ',', 'while', 'the', 'Americans', 'more', 'readily', 'kept', 'such', 'use', 'until', 'more', 'recently', ',', 'and', 'still', 'maintain', 'it', 'more', 'than', 'Britons', '.', 'The', 'classic', 'example', ',', 'considered', 'by', 'their', 'American', 'counterparts', 'quite', 'curious', ',', 'was', 'the', 'maintenance', 'of', 'the', 'internal', 'comma', 'in', 'a', 'British', 'organisation', 'of', 'secret', 'agents', 'called', 'the', 'Special', 'Operations', ',', 'Executive', 'S.O.', 'E.', 'which', 'is', 'not', 'found', 'in', 'histories', 'written', 'after', 'about', '1960', 'But', 'before', 'that', ',', 'many', 'Britons', 'were', 'more', 'scrupulous', 'at', 'maintaining', 'the', 'French', 'form', '.', 'In', 'French', ',', 'the', 'period', 'only', 'follows', 'an', 'abbreviation', 'if', 'the', 'last', 'letter', 'in', 'the', 'abbreviation', 'is', 'not', 'the', 'last', 'letter', 'of', 'its', 'antecedent', '.', 'M.', 'is', 'the', 'abbreviation', 'for', 'monsieur', 'while', 'Mme', 'is', 'that', 'for', 'madame', '.', 'Like', 'many', 'other', 'cross-channel', 'linguistic',

'acquisitions', ',', 'many', 'Britons', 'readily', 'took', 'this', 'up', 'and', 'followed', 'this', 'rule', 'themselves', ',', 'while', 'the', 'Americans', 'took', 'a', 'simpler', 'rule', 'and', 'applied', 'it', 'rigorously', '., 'Over', 'the', 'years', '., 'however', '., 'the', 'lack', 'of', 'convention', 'in', 'some', 'style', 'guides', 'has', 'made', 'it', 'difficult', 'to', 'determine', 'which', 'two-word', 'abbreviations', 'should', 'be', 'abbreviated', 'with', 'periods', 'and', 'which', 'should', 'not', '., 'The', 'U.S.', 'media', 'tend', 'to', 'use', 'periods', 'in', 'two-word', 'abbreviations', 'like', 'United', 'States', '(', 'U.S.', ')', '., 'but', 'not', 'personal', 'computer', '(', 'PC', ')', 'or', 'television', '(', 'TV', ')', '., 'Many', 'British', 'publications', 'have', 'gradually', 'done', 'away', 'with', 'the', 'use', 'of', 'periods', 'in', 'abbreviations', '., 'Minimization', 'of', 'punctuation', 'in', 'typewritten', 'material', 'became', 'economically', 'desirable', 'in', 'the', '1960s', 'and', '1970s', 'for', 'the', 'many', 'users', 'of', 'carbon-film', 'ribbons', 'since', 'a', 'period', 'or', 'comma', 'consumed', 'the', 'same', 'length', 'of', 'non-reusable', 'expensive', 'ribbon', 'as', 'did', 'a', 'capital', 'letter', '., 'Widespread', 'use', 'of', 'electronic', 'communication', 'through', 'mobile', 'phones', 'and', 'the', 'Internet', 'during', 'the', '1990s', 'allowed', 'for', 'a', 'marked', 'rise', 'in', 'colloquial', 'abbreviation', '., 'This', 'was', 'due', 'largely', 'to', 'increasing', 'popularity', 'of', 'textual', 'communication', 'services', 'such', 'as', 'instant', '-', 'and', 'text', 'messaging', '., 'SMS', '., 'for', 'instance', '., 'supports', 'message', 'lengths', 'of', '160', 'characters', 'at', 'most', '(', 'using', 'the', 'GSM', '03', '., '38', 'character', 'set', ')', '., 'This', 'brevity', 'gave', 'rise', 'to', 'an', 'informal', 'abbreviation', 'scheme', 'sometimes', 'called', 'Textese', '., 'with', 'which', '10', 'or', 'more', 'of', 'the', 'words', 'in', 'a', 'typical', 'SMS', 'message', 'are', 'abbreviated', '., 'More', 'recently', 'Twitter', '., 'a', 'popular', 'social', 'networking', 'service', '., 'began', 'driving', 'abbreviation', 'use', 'with', '140', 'character', 'message', 'limits

2. با استفاده از عبارات منظم و الگوریتم Porter (الگوریتم کاملتر را از اینترنت استخراج کنید) ریشه کلمات (Stemming) را بدست آورید.

جهت نوشتن قوانین الگوریتم porter از فایل ضمیمه شده برای پیدا کردن قوانین استفاده شده. تابع `apply_porter_stemming` جهت چک کردن مرحله به مرحله تکست و عمل `stemming` بر روی متن مورد نظر. در قانون اول که شامل سه بخش میشود، در بخش اول کلمات جمع به ساده و در بخش دوم افعال به شکل ساده تغییر کردند و در بخش سوم حرف `y` اضافه در انتهای فعل را به `a` تغییر داده شده. در قانون دوم بیشتر `suffix`ها را به کلمات ساده تبدیل شدند و عوض شدند در قانون سوم نیز به همین ترتیب یا به کلمات ساده تر تبدیل شدند یا حرف اضافه از آنها حذف شد. در قانون چهارم حرف اضافی `ly` که به عنوان متمم گاها در کنار کلمات قرار میگیرد حذف شده.

- جهت به کار گیری تابع نوشته شده، ابتدا توکن های ساخته شده در مرحله قبل را به تکست تبدیل شدند و بعد تابع را پیاده سازی کرده و پرینت میگیریم.

- رفرنس: <https://people.scs.carleton.ca/~armyunis/projects/KAPI/porter.pdf>

```
# Join tokens to form the text for stemming
text = " ".join(tokens)

# Apply stemming
stemmed_text = apply_porter_stemming(text)

print(stemmed_text)
```

assignment1

```
def apply_porter_stemming(text: Any) -> str :
```

- Stemming function

```

11 def apply_porter_stemming(text):
12     # Rule 1a: Plurals
13     text = re.sub(r'(sses)\b', r'ss', text)
14     text = re.sub(r'(ies)\b', r'i', text)
15     text = re.sub(r'([^s])s\b', r'\1', text)
16
17     # Rule 1b: Past tense and gerund forms
18     text = re.sub(r'(ed|ing)\b', '', text)
19
20     # Rule 1c: Change 'y' to 'i' when preceded by a vowel
21     text = re.sub(r'([aeiou])y\b', r'\1i', text)
22
23     # Rule 2: Double suffixes
24     text = re.sub(r'ization\b', 'ize', text)
25     text = re.sub(r'ational\b', 'ate', text)
26     text = re.sub(r'ation\b', 'ate', text)
27     text = re.sub(r'ator\b', 'ate', text)
28     text = re.sub(r'alism\b', 'al', text)
29     text = re.sub(r'iveness\b', 'ive', text)
30     text = re.sub(r'fulness\b', 'ful', text)
31     text = re.sub(r'ousness\b', 'ous', text)
32     text = re.sub(r'aliti\b', 'al', text)
33     text = re.sub(r'iviti\b', 'ive', text)
34     text = re.sub(r'biliti\b', 'ble', text)
35
36     # Rule 3:
37     text = re.sub(r'icate\b', 'ic', text)
38     text = re.sub(r'ative\b', '', text)
39     text = re.sub(r'alize\b', 'al', text)
40     text = re.sub(r'iciti\b', 'ic', text)
41
42     text = re.sub(r'ical\b', 'ic', text)
43     text = re.sub(r'ful\b', '', text)
44     text = re.sub(r'ness\b', '', text)
45
46     # Rule 4: adverb suffix -ly
47     text = re.sub(r'ly\b', '', text)
48
49     return text

```

3. با استفاده از عبارات منظم متن را Nomalize کنید.

برای normalization چندین مرحله از روی رفرنس ذکر شده انجام شد.

- Lowercase
- Punctuation removal
- Remove extra whitespace
- Common contractions
- stop words (and, the, is, in, of, are)

```

50 def normalization(text):
51     # Step1: Lowercase
52     text = text.lower()
53
54     # Step2: Remove punctuation
55     text = re.sub(r'^\w\s.-', '', text)
56
57     # Step3: Remove extra whitespace
58     text = text = re.sub(r'\s+', ' ', text).strip()
59
60     # Step4: Handle common contractions
61     contractions = {
62         "don't": "do not",
63         "doesn't": "does not",
64         "can't": "cannot",
65         "won't": "will not",
66         "i'm": "i am",
67         "you're": "you are",
68         "we're": "we are",
69         "they're": "they are",
70         "isn't": "is not",
71         "aren't": "are not",
72         "wasn't": "was not",
73         "weren't": "were not",
74         "it's": "it is",
75         "that's": "that is",
76         "there's": "there is",
77     }
78
79     for contraction, expanded_form in contractions.items():
80         text = re.sub(r'\b' + contraction + r'\b', expanded_form, text)

```

```

81
82     # Step5: Remove stop words
83     tokens = text.split()
84     stop_words = set(["and", "the", "is", "in", "of", "on", "an", "a"])
85     tokens = [word for word in tokens if word not in stop_words]
86
87     return text
88

```

در نهایت متن به شکل زیر قرار گرفت:

after world war ii the british great reduc the use of the full stop and other punctuate point after abbreviate in at least semi-formal writ while the american more readi kept such use until more recent and still maintain it more than briton . the classic example consider by their american counterpart quite curiou wa the maintenance of the internal comma in a british organisate of secret agent call the special operate executive s.o. e which i not found in histori written after about 1960 . but before that many briton were more scrupulou at maintain the french form . in french the period on follow an abbreviate if the last letter in the abbreviate i not the last letter of it antecedent m. i the abbreviate for monsieur while mme i that for madame . like many other cross-channel linguistic acquisition many briton readi took thi up and follow thi rule themself while the american took a simpler rule and appli it rigorous . over the year however the lack of convention in some style guide ha made it difficult to determine which two-word abbreviate should be abbreviat with period and which should not . the u.s. media tend to use period in two-word abbreviate like unit state u.s. but not personal computer pc or television tv . many british public have gradual done awai with the use of period in abbreviate . minimize of punctuate in typewritten material became economical desirable in the 1960 and 1970 for the many user of carbon-film ribbon since a period or comma consum the same length of non-reusable expensive ribbon a did a capital letter . widespread use of electronic communic through mobile phone and the internet dur the 1990 allow for a mark rise in colloquial abbreviate . thi wa due large to increas popularity of textual communic service such a instant - and text messag . sms for instance support message length of 160 character at most us the gsm 03 . 38 character set . thi brevity gave rise to an informal abbreviate scheme sometime call textese with which 10 or more of the word in a typic sms message are abbreviat . more recent twitter a popular social network service began driv abbreviate use with 140 character message limit .

در این متن چندین ایراد بزرگ وجود دارد که به درستی کلمات ساده سازی نشده اند.

- چندین s، ed که به جهت تغییر اسم جمع حذف قرار بود بشوند باعث شدند کل کلمه اشتباه در نظر گرفته بشود.
- در lowercase کردن تمامی کلمات باعث میشود گاهی برخی از اسامی اشتباه استنباط شوند.
- کلماتی که دارای – بودند معنای خود را از دست دادند.
- تبدیل کردن ing باعث شده است که کلمه کامل از بین برود.