

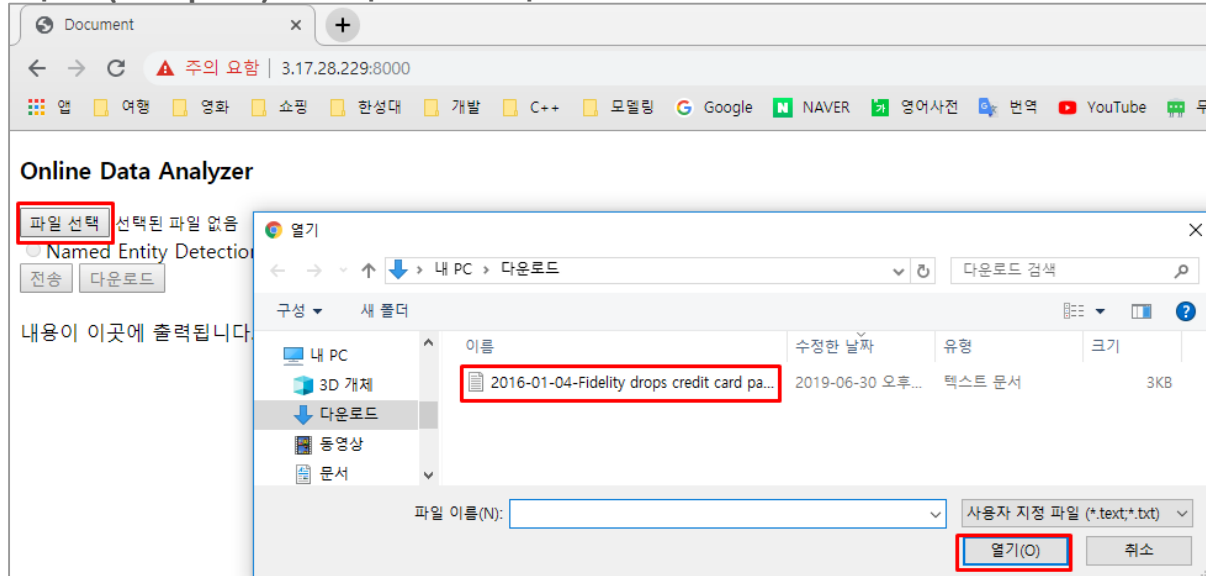


Online Data Analyzer 사용 방법

Jin Won,
School of Computer Engineering,
Hansung University

Contents

- 웹 페이지 접속
- 파일(corpus) 선택 -> 열기



Named Entity Detection

Online Data Analyzer

파일 선택 2016-01-04-Fi...f America.txt

☒ Named Entity Detection ☐ N-Grams Extraction ☐ Word Pair Extraction ☐ Phrase Extraction

```
1 from application.nlp import comprehend
2 from nltk import sent_tokenize
3
4 def run(article):
5     # 기사를 라인 별로 분할하여 리스트로 저장한다.
6     sentences = sent_tokenize(article)
7
8     # Map 객체로써 key는 entity 명, value는 타입이다.
9     # ex. key: NEW YORK, value: LOCATION
10    dictMap = {}
11
12    # 결과를 저장하는 string이다.
13    retVal = ''
14
15    # sentences 내에 있는 모든 sentence에 대해
16    for sentence in sentences:
17        # entity를 조사한다.
18        entities = comprehend.detect_entities(Text=sentence, LanguageCode='en')
19
20        # entities 내에 있는 모든 entity에 대해
21        for entity in entities["Entities"]:
22            # 이름과 type을 저장한다.
23            key = entity["Text"]
24            value = entity["Type"]
25
26            # dictMap에 해당 이름의 entity가 존재하지 않는 경우에 한해서
27            # entity를 추가하고, retVal에 반영한다.
28            if key not in dictMap:
29                dictMap[key] = value
30                row = (key + ': ' + value + '\n\r\n')
31                retVal += row
32
33    return retVal
```

Named Entity Detection

• input

BOSTON (Reuters) - Fidelity Investments said on Monday it is dropping long-time credit card partners American Express Co and Bank of America Corp, ending a 12-year partnership that has generated billions of dollars in fees. Boston-based Fidelity, which has 24 million customers, said its new partners will be U.S. Bancorp and Visa Inc, effective Monday. The exclusive alliance will provide Visa branded credit-card products to U.S. consumers, including Fidelity customers. The switch is another setback for American Express, already reeling from its lost deal with warehouse club retailer Costco Wholesale Corp. AmEx said earlier this year the loss of the Costco contract would hurt profit for two years. AmEx shares are off 25 percent over the past year and were down 2.9 percent Monday morning. Many other major U.S. financial stocks were also lower, with shares of Bank of America, Visa Inc and U.S. Bancorp all down by more than 2 percent. Ram Subramaniam, president of Fidelity's retail brokerage business, did not give any specific reason for ending the partnership with American Express and Bank of America. "It's been a long, good partnership," he said. A spokeswoman for American Express said the Fidelity portfolio accounts for less than 1 percent of billings. A Bank of America spokeswoman said the agreement not to continue the relationship with American Express was a mutual decision between the two companies. "Over the past several years, Bank of America has been exiting from our financial institutions card business where Bank of America has limited opportunity to deepen customer relationships, and this move is consistent with that strategy," she wrote via email. Since 2003, Fidelity has offered 2 percent cash back credit cards with American Express and Bank of America's FIA Card Services. During that time Fidelity customers have earned \$1.1 billion cash rewards. The new alliance will feature cards with chip security technology, with access to digital wallets that include Apple Pay, Samsung Pay and Android pay. The new card program will issue the Fidelity Rewards Visa Signature Card and the Fidelity Investments 529 College Rewards Visa Signature Card, where card members can earn unlimited 2 percent cash back with no annual fees, caps or categories when directing rewards into eligible Fidelity accounts. U.S. Bank also has agreed to acquire Fidelity's existing co-brand credit card portfolio with about \$1.7 billion in associated balances. Additional reporting by Richa Naidu in Bangalore and Dan Freed in New York; Editing by Alan Crosby and David Gregorio

output

BOSTON: LOCATION
Reuters: ORGANIZATION
Fidelity Investments: ORGANIZATION
Monday: DATE
American Express Co: ORGANIZATION
Bank of America Corp: ORGANIZATION
12-year: QUANTITY
billions of dollars: QUANTITY
Boston: LOCATION
Fidelity: ORGANIZATION
24 million customers: QUANTITY
U.S.: LOCATION
Bancorp: ORGANIZATION
Visa Inc: ORGANIZATION
Visa: ORGANIZATION
U.S.: LOCATION
American Express: ORGANIZATION
Costco Wholesale Corp.: ORGANIZATION
AmEx: ORGANIZATION
earlier this year: DATE
Costco: ORGANIZATION
two years: QUANTITY
25 percent: QUANTITY
past year: DATE
2.9 percent: QUANTITY
Monday morning: DATE
Bank of America: ORGANIZATION
more than 2 percent: QUANTITY
Ram Subramaniam: PERSON
less than 1 percent: QUANTITY
two companies: QUANTITY
past several years: DATE
2003: DATE

N-Grams Extraction

Online Data Analyzer

파일 선택 2016-01-04-Fi...f America.txt

☐ Named Entity Detection ☒ N-Grams Extraction ☐ Word Pair Extraction ☐ Phrase Extraction

N-Grams Extraction Parameters

토큰 수: 2

빈도수 임계값: 3

전송

다운로드

```
1 from nltk import *
2
3 # line 단위의 sentence를 입력받아 빈도 분포 객체를 반환한다.
4 # @param sentence: 빈도 분포를 분석할 문장
5 # @param numTokens: 빈도 분포를 생성할 기준 토큰 수
6 def getFreqDist(sentence, numTokens):
7     tokenList = word_tokenize(sentence)
8     ngramList = ngrams(tokenList, numTokens)
9
10    return FreqDist(ngramList)
11
12 def run(article, numTokens, freqThreshold):
13     # 텍스트 파일을 라인 별로 분할하여 리스트로 저장
14     sentences = article.splitlines()
15
16     # sentences의 모든 빈도 분포를 분석한 결과가 저장된다.
17     # 이 객체는 Map 객체이며, key는 복합어, value는 빈도 수이다.
18     globalFreqDist = {}
19
20     # sentences 내에 있는 모든 sentence에 대하여
21     for sentence in sentences:
22         # 현재 sentence에 대한 빈도 분포 분석
23         localDist = getFreqDist(sentence, numTokens)
24
25         # localDist에 존재하는 빈도 분포를 조사
26         # key: 복합어 / value: 빈도 수
27         for key, value in localDist.items():
28
29             # 만약 globalFreqDist에 이미 존재하는 복합어라면
30             if key in globalFreqDist:
31                 # 해당 복합어의 빈도 수에 value를 누적
32                 globalFreqDist[key] = globalFreqDist[key] + value
33             # 처음 나타난 복합어라면
34             else:
35                 # globalFreqDist 맵에 새로 추가
36                 globalFreqDist[key] = value
37
38     retVal = ''
```

N-Grams Extraction

• input

BOSTON (Reuters) - Fidelity Investments said on Monday it is dropping long-time credit card partners American Express Co and Bank of America Corp, ending a 12-year partnership that has generated billions of dollars in fees. Boston-based Fidelity, which has 24 million customers, said its new partners will be U.S.

Bancorp and Visa Inc, effective Monday.

The exclusive alliance will provide Visa branded credit-card products to U.S. consumers, including Fidelity customers.

The switch is another setback for American Express, already reeling from its lost deal with warehouse club retailer Costco Wholesale Corp.

AmEx said earlier this year the loss of the Costco contract would hurt profit for two years.

AmEx shares are off 25 percent over the past year and were down 2.9 percent Monday morning.

Many other major U.S. financial stocks were also lower, with shares of Bank of America, Visa Inc and U.S.

Bancorp all down by more than 2 percent.

Ram Subramaniam, president of Fidelity's retail brokerage business, did not give any specific reason for ending the partnership with American Express and Bank of America.

"It's been a long, good partnership," he said.

A spokeswoman for American Express said the Fidelity portfolio accounts for less than 1 percent of billings.

A Bank of America spokeswoman said the agreement not to continue the relationship with American Express was a mutual decision between the two companies.

"Over the past several years, Bank of America has been exiting from our financial institutions card business where Bank of America has limited opportunity to deepen customer relationships, and this move is consistent with that strategy," she wrote via email.

Since 2003, Fidelity has offered 2 percent cash back credit cards with American Express and Bank of America's FIA Card Services. During that time Fidelity customers have earned \$1.1 billion cash rewards.

The new alliance will feature cards with chip security technology, with access to digital wallets that include Apple Pay, Samsung Pay and Android pay.

The new card program will issue the Fidelity Rewards Visa Signature Card and the Fidelity Investments 529 College Rewards Visa Signature Card, where card members can earn unlimited 2 percent cash back with no annual fees, caps or categories when directing rewards into eligible Fidelity accounts.

U.S.

Bank also has agreed to acquire Fidelity's existing co-brand credit card portfolio with about \$1.7 billion in associated balances.

Additional reporting by Richa Naidu in Bangalore and Dan Freed in New York; Editing by Alan Crosby and David Gregorio

output

('American', 'Express'), 6
('and', 'Bank'), 3
('Bank', 'of'), 7
('of', 'America'), 7
('U.S', '.'), 3
('2', 'percent'), 3
('with', 'American'), 3
('the', 'Fidelity'), 3

Word Pair Extraction

code

```
1 from nltk import sent_tokenize
2 from nltk import RegexpParser
3 from nltk import word_tokenize
4 from nltk import pos_tag
5 from nltk import tree
6 from re import search
7
8 def run(article, pos1, pos2):
9     # 텍스트 파일을 문장 단위로 분할하여 sentenceList로 저장
10    sentenceList = sent_tokenize(article)
11
12    # 재정의의 품사 패턴
13    noun = 'FW|NN' # 명사
14    verb = 'MD|VB' # 동사
15    adj = 'JJ|VBN' # 형용사
16    adv = 'W?RB|*RP' # 부사
17    det = 'DT?' # 한정사
18    pron = 'EX|PRP|*WP' # 대명사
19    prep = '^TO|*IN' # 전치사
20    conj = 'CC' # 접속사
21    int_ = 'UH' # 감탄사
22
23    # word pair 패턴 정의
24    patternList = ('pattern1: {<%s><%s>}' % (pos1, pos2))
25
26    # 결과를 저장하는 string이다.
27    retVal = ''
28
29    # 파서 생성
30    parser = RegexpParser(patternList)
31
```

```
32 # 리스트 내 모든 문장에 대하여
33 for sentence in sentenceList:
34
35     # 공백을 기준으로 tokenizing
36     tokenList = word_tokenize(sentence)
37
38     # 리스트 내 모든 토큰들을 품사별로 매핑
39     taggedList = pos_tag(tokenList)
40
41     # (토큰, 재정의된 품사) 튜플을 저장할 리스트
42     retaggedList = []
43
44     # POS 재정의
45     for tagged in taggedList:
46
47         if search(noun, tagged[1]):
48             retaggedList.append((tagged[0], 'NOUN'))
49         elif search(verb, tagged[1]):
50             retaggedList.append((tagged[0], 'VERB'))
51         elif search(adj, tagged[1]):
52             retaggedList.append((tagged[0], 'ADJ'))
53         elif search(adv, tagged[1]):
54             retaggedList.append((tagged[0], 'ADV'))
55         elif search(det, tagged[1]):
56             retaggedList.append((tagged[0], 'DET'))
57         elif search(pron, tagged[1]):
58             retaggedList.append((tagged[0], 'PRON'))
59         elif search(pre, tagged[1]):
60             retaggedList.append((tagged[0], 'PREP'))
61         elif search(conj, tagged[1]):
62             retaggedList.append((tagged[0], 'CONJ'))
63         elif search(int_, tagged[1]):
64             retaggedList.append((tagged[0], 'INT'))
65         else:
66             retaggedList.append((tagged[0], tagged[1]))

```

Online Data Analyzer

파일 선택 2016-01-04-Fi...f America.txt

Named Entity Detection N-Grams Extraction Word Pair Extraction Phrase Extraction

Word Pair Extraction Parameters 명사(NOUN) 동사(VERB)

전송 다운로드 명사(NOUN) 동사(VERB)

명용사(ADJ) 부사(ADV)

한정사(DET) 대명사(PRON)

전치사(PREP) 접속사(CONJ)

감탄사(INT)

BOSTON (Reuters) - Fidelity Investments said on

Bank of America Corp, ending a 12-year partners

Boston-based Fidelity, which has 24 million custo

Bancorp and Visa Inc, effective Monday.

The exclusive alliance will provide Visa branded c

The switch is another setback for American Expre

opping long-time credit card p
nerated billions of dollars in fe
ew partners will be U.S.

ucts to U.S. consumers, includi
ng from its last deal with w

```
68 # 파싱
69 parsedTree = parser.parse(retaggedList)
70
71 # 파일에 demo1 출력
72 for subtree in parsedTree:
73     if isinstance(subtree, tree.Tree):
74         retVal += (subtree[0][0] + ' [' + subtree[0][1] + ' ] ' )
75         retVal += (subtree[1][0] + ' [' + subtree[1][1] + ' ]\n')
76
77 return retVal

```

Word Pair Extraction

• input

BOSTON (Reuters) - Fidelity Investments said on Monday it is dropping long-time credit card partners American Express Co and Bank of America Corp, ending a 12-year partnership that has generated billions of dollars in fees. Boston-based Fidelity, which has 24 million customers, said its new partners will be U.S. Bancorp and Visa Inc, effective Monday. The exclusive alliance will provide Visa branded credit-card products to U.S. consumers, including Fidelity customers. The switch is another setback for American Express, already reeling from its lost deal with warehouse club retailer Costco Wholesale Corp. AmEx said earlier this year the loss of the Costco contract would hurt profit for two years. AmEx shares are off 25 percent over the past year and were down 2.9 percent Monday morning. Many other major U.S. financial stocks were also lower, with shares of Bank of America, Visa Inc and U.S. Bancorp all down by more than 2 percent. Ram Subramaniam, president of Fidelity's retail brokerage business, did not give any specific reason for ending the partnership with American Express and Bank of America. "It's been a long, good partnership," he said. A spokeswoman for American Express said the Fidelity portfolio accounts for less than 1 percent of billings. A Bank of America spokeswoman said the agreement not to continue the relationship with American Express was a mutual decision between the two companies. "Over the past several years, Bank of America has been exiting from our financial institutions card business where Bank of America has limited opportunity to deepen customer relationships, and this move is consistent with that strategy," she wrote via email. Since 2003, Fidelity has offered 2 percent cash back credit cards with American Express and Bank of America's FIA Card Services. During that time Fidelity customers have earned \$1.1 billion cash rewards. The new alliance will feature cards with chip security technology, with access to digital wallets that include Apple Pay, Samsung Pay and Android pay. The new card program will issue the Fidelity Rewards Visa Signature Card and the Fidelity Investments 529 College Rewards Visa Signature Card, where card members can earn unlimited 2 percent cash back with no annual fees, caps or categories when directing rewards into eligible Fidelity accounts. U.S. Bank also has agreed to acquire Fidelity's existing co-brand credit card portfolio with about \$1.7 billion in associated balances. Additional reporting by Richa Naidu in Bangalore and Dan Freed in New York; Editing by Alan Crosby and David Gregorio

output

Investments [NOUN] said [VERB]
partners [NOUN] will [VERB]
alliance [NOUN] will [VERB]
Visa [NOUN] branded [VERB]
switch [NOUN] is [VERB]
AmEx [NOUN] said [VERB]
contract [NOUN] would [VERB]
shares [NOUN] are [VERB]
stocks [NOUN] were [VERB]
Express [NOUN] said [VERB]
spokeswoman [NOUN] said [VERB]
Express [NOUN] was [VERB]
America [NOUN] has [VERB]
America [NOUN] has [VERB]
move [NOUN] is [VERB]
Fidelity [NOUN] has [VERB]
customers [NOUN] have [VERB]
alliance [NOUN] will [VERB]
program [NOUN] will [VERB]
members [NOUN] can [VERB]

Phrase Extraction

Online Data Analyzer

파일 선택 2016-01-04-Fi...f America.txt

☐ Named Entity Detection ☐ N-Grams Extraction ☐ Word Pair Extraction ☒ Phrase Extraction

전송

다운로드

*현재는 명사구만 추출되며
chunk 정규식을 어떻게 입력 받을지에 따라
개선할 수 있습니다.

```
1 from nltk import word_tokenize
2 from nltk import pos_tag
3 from nltk import RegexpParser
4 from nltk import tree
5
6 def run(article):
7     # 컨텐츠 내에 있는 모든 단어들을 공백 문자열 기준으로 Tokenizing
8     rawTokens = word_tokenize(article)
9
10    # Part of Speech(POS) tagging
11    # rawTokens(list)내에 저장된 모든 Token들을 POS 별로 매핑
12    taggedTokens = pos_tag(rawTokens)
13
14    # Chunk 정규식 정의
15    # Chunk 정규식은 절대적인 기준이 없으며, 개발자가 직접 원문의 내용을 분석하고
16    # 상황에 따라 별도로 정의해주어야 한다.
17    regexes = """
18        NP: {<DT|PP\$>?<JJ>*<NN.*>+} # noun phrase
19        PP: {<IN><NP>} # prepositional phrase
20        VP: {<MD>?<VB.*><NP|PP>} # verb phrase
21    """
22
23    # 파서 생성
24    parser = RegexpParser(regexes)
25
26    # 파싱
27    parsedTree = parser.parse(taggedTokens)
28
29    # 결과를 저장하는 string이다.
30    retVal = ''
31
32    for subtree in parsedTree:
33        if isinstance(subtree, tree.Tree) and subtree.label() == "NP":
34            retVal += (str(subtree) + '\r\n')
35
36    return retVal
```

Phrase Extraction

• input

BOSTON (Reuters) - Fidelity Investments said on Monday it is dropping long-time credit card partners American Express Co and Bank of America Corp, ending a 12-year partnership that has generated billions of dollars in fees. Boston-based Fidelity, which has 24 million customers, said its new partners will be U.S. Bancorp and Visa Inc, effective Monday. The exclusive alliance will provide Visa branded credit-card products to U.S. consumers, including Fidelity customers. The switch is another setback for American Express, already reeling from its lost deal with warehouse club retailer Costco Wholesale Corp. AmEx said earlier this year the loss of the Costco contract would hurt profit for two years. AmEx shares are off 25 percent over the past year and were down 2.9 percent Monday morning. Many other major U.S. financial stocks were also lower, with shares of Bank of America, Visa Inc and U.S. Bancorp all down by more than 2 percent. Ram Subramaniam, president of Fidelity's retail brokerage business, did not give any specific reason for ending the partnership with American Express and Bank of America. "It's been a long, good partnership," he said. A spokeswoman for American Express said the Fidelity portfolio accounts for less than 1 percent of billings. A Bank of America spokeswoman said the agreement not to continue the relationship with American Express was a mutual decision between the two companies. "Over the past several years, Bank of America has been exiting from our financial institutions card business where Bank of America has limited opportunity to deepen customer relationships, and this move is consistent with that strategy," she wrote via email. Since 2003, Fidelity has offered 2 percent cash back credit cards with American Express and Bank of America's FIA Card Services. During that time Fidelity customers have earned \$1.1 billion cash rewards. The new alliance will feature cards with chip security technology, with access to digital wallets that include Apple Pay, Samsung Pay and Android pay. The new card program will issue the Fidelity Rewards Visa Signature Card and the Fidelity Investments 529 College Rewards Visa Signature Card, where card members can earn unlimited 2 percent cash back with no annual fees, caps or categories when directing rewards into eligible Fidelity accounts. U.S. Bank also has agreed to acquire Fidelity's existing co-brand credit card portfolio with about \$1.7 billion in associated balances. Additional reporting by Richa Naidu in Bangalore and Dan Freed in New York; Editing by Alan Crosby and David Gregorio

output

(NP BOSTON/NNP)
(NP Reuters/NNP)
(NP Fidelity/NN Investments/NNPS)
(NP Bank/NNP)
(NP Boston-based/JJ Fidelity/NNP)
(NP customers/NN)
(NP new/JJ partners/NN)
(NP Visa/NNP Inc/NNP)
(NP effective/JJ Monday/NNP)
(NP The/DT exclusive/JJ alliance/NN)
(NP U.S./NNP consumers/NN)
(NP The/DT switch/NN)
(NP this/DT year/NN)
(NP the/DT loss/NN)
(NP years/NN)
(NP AmEx/NNP shares/NN)
(NP percent/NN)
(NP percent/NN Monday/NNP morning/NN)
(NP other/JJ major/JJ U.S./NNP)
(NP financial/JJ stocks/NN)
(NP Visa/NNP Inc/NNP)
(NP U.S./NNP Bancorp/NNP)
(NP percent/NN)
(NP Ram/NNP Subramaniam/NNP)
(NP president/NN)
(NP retail/JJ brokerage/NN business/NN)
(NP Bank/NNP)
(NP good/JJ partnership/NN)