

[R] R을 통한 텍스트마이닝에서 워드클라우드 까지 — 나무늘보의 개발 블로그

노트북: blog

만든 날짜: 2020-10-06 오후 4:44

URL: <https://continuous-development.tistory.com/53?category=793392>

R

[R] R을 통한 텍스트마이닝에서 워드클라우드 까지

2020. 8. 5. 23:55 수정 삭제 공개

#텍스트 마이닝

시작전에 자바 jdk를 깔고 환경변수를 설정해줘야 한다. jdk를 통해 KoNLP 통해 돌리기 때문이다. 이부분이 선행되어야 나머지가 된다.

```
9
10 install.packages(c("hash", "tau", "Sejong",
11                   "RSQLite", "devtools", "bit",
12                   "rex", "lazyeval", "htmlwidgets",
13                   "crosstalk", "promises", "later",
14                   "sessioninfo", "xopen", "bit64",
15                   "blob", "DBI", "memoise", "plogr",
16                   "covr", "DT", "rcmdcheck", "rversions"),
17                   type = "binary")
18
19 # github에서 패키지 설치
20 install.packages("remotes")
21 # 64bit 에서만 동작합니다.
22 remotes::install_github('haven-jeon/KoNLP',
23                        upgrade = "never",
24                        INSTALL_opts=c("--no-multiarch")) #konlp는 직접인스톨하면 에러가 나서 이런 방식으로 인스톨해야 된다.
25
26 devtools::install_github("r-pkgs/usethis")
```

KoNLP는 일반적인 인스톨을 해서는 안되고 위와같이 github라는 것을 통해서 받아야 한다.

```

22
23 #기본 라이브러리 로딩
24 Sys.setenv(JAVA_HOME='C:\\Program Files\\JAVA\\jdk1.8.0_121')
25 library(rJava)
26 library(KoNLP)
27
28 install.packages("tm")
29 install.packages("wordcloud")
30 library(tm) # 전처리 용도 라이브러리
31 library(wordcloud)
32 library(RColorBrewer) # 팔레트 역할 (색깔)
33
34

```

없는부분은 install하자 java home의 위치나 버전은 사람마다 다를 수도 있으니 해당 경로를 맞춰주자.

```

35 # 명사 , 형용사등을 추출한 사진
36 ?useSejongDic()
37
38 useSejongDic(text)
39
40 text <-|"최근 이슈가 되고 있는 빅데이터에 대한 이해와 활용을 위해 데이터 과학(Data Science)의 측면에서 접근한다.
41 빅데이터는 통계학을 비롯한 경영, IT 등의 다양한 분야들이 서로 결합되어 있고 그 정의가 다양하지만, 본 강의는 데이터 분석을
42 기반으로 하는 과학적 의사결정의 관점에서 바라보고자 한다. 빅데이터에 대한 이해를 위해 실제 사례들을 살펴보고,
43 데이터를 통해 의사결정에 유용한 정보 및 지식을 찾는 과정을 이해한다. 나아가 빅데이터 분석에서 필수적으로 언급되고 있는
44 R 통계프로그램을 소개하고 이를 분석에 활용할 수 있게 한다."
45
46 nouns <- extractNoun(text)
47

```

test라는 변수에 문자열을 넣어준다.

extractNoun을 통해 명사를 추출해낸다.

```

51
52 nouns <- extractNoun(text) #명사를 뽑아내는 기능
53
54

```

```

> nouns <- extractNoun(text)
> nouns
[1] "이슈"      "빅데이터에" "이해"      "활용"      "데이터"    "과학(Data" "Science"   "측면"      "검급"      "빅데이터는" "통계학"    "비롯"      "한"      "경영"
[15] "IT"        "등"         "다양"      "문"        "분야"      "들"        "결합"      "다양"      "강의"      "데이터"    "분석"      "기반"      "과학"
[29] "적"        "의사결정"   "관점"      "빅데이터에" "이해"      "실제"      "사례"      "들"        "의사결정" "유용"      "현"        "정보"    "지식"
[43] "과학"     "이해"      "빅데이터"  "분석"      "일수"      "적"        "언급"      "R"         "통계"      "프로그램"  "소개"      "이"      "분석"
[57] "수"

```

이제 뽑아낸 명사를 전처리 하는 과정을 거친다.

#전처리

명사를 추출한 데이터에서 nchar()를 통해 두글자 이상의 글자만 가져온다. 한글자로 만든 것중에 등 한 이런것들을 빼기 위해서이다.

```
54
55 # nchar() - 문자열 길이 함수
56 nouns <- nouns[nchar(nouns) >= 2]
57 nouns
58
59
```

table 함수를 사용해서 각각의 데이터가 몇번씩 나왔는지 확인한다.

```
> # 빈도표
> wFreq <- table(nouns)
> wFreq
nouns
IT      Science   강의   결합   경영   과정   과학   과학(Data   관점   기반   다양   데이터   들이   분석   분야   비롯   빅데이터   빅데이터는
빅데이터 1       1       1       1       1       1       1       1       1       1       1       1       1       1       1       1       1
2       사례 1       소개 1       실제 1       언급 1       유용 1       의사결정 2       이슈 1       이해 3       접근 1       정보 1       정의 1       지식 1       측면 1       통계 1       통계학 1       프로그램 1       필수 1       활용 2
```

빅데이터~ 는 명사로 뽑으면 빅데이터 단어가 끝이기 때문에 전처리가 제대로 안된 부분에 대해서는 이런식으로 전처리를 해준다.

```
58
59 #gsub은 바꾸는 기능을 한다. .은 한글자 그리고 *를 쓰면 무한대를 말한다. 그래서 빅데이터/는 이런 것들을 빅데이터 로 바꿔준다.
60 nouns <- gsub("빅데이터.*", "빅데이터", nouns) #
61 nouns
62
63 wFreq <- table(nouns)
64 wFreq
65
```

```
nouns
IT      Science   강의   결합   경영   과정   과학   과학(Data   관점   기반
1       1       1       1       1       1       1       1       1       1
다양   데이터   들이   분석   분야   비롯   빅데이터   사례   소개   실제
2       3       1       3       1       1       4       1       1       1
언급   유용   의사결정   이슈   이해   접근   정보   정의   지식   측면
1       1       2       1       3       1       1       1       1       1
통계   통계학   프로그램   필수   활용
1       1       1       1       2
```

테이블 데이터를 WFreq로 확인한다.

이제 테이블 데이터를 가지고 워드클라우드를 만든다.

```

72
73 names(wFreq)
74
75 pal <- brewer.pal(6,"Accent") #6가지 Accent 를 나타낸다.
76
77 wordcloud(words = names(wFreq), # 빈도를 가지고 있는 속성값들을 넣어준다.
78           freq = wFreq,         #빈도를 가지고 출력할 데이터를 넣는다.
79           min.freq = 1,         #최소 빈도수를 정한다. 여기서 1이면 1개 이상은 모두 출력된다.
80           random.order = F,     #임의의 순서로 작성하는데 F일때는 감소 빈도이다
81           colors=pal)          #색상을 넣는다.
82 ?wordcloud
83

```



예제

기본적인 텍스트 마이닝 단계는 text -> corpus(말뭉치) -> TDM(단어에 따른 매트릭스) -> 매트릭스로 형변환 -> 워드클라우드 로 나뉜다. 전체적인 로직은 이렇다.

```

159 #데이터 읽어오기
160 data <- readLines(file.choose(),encoding = 'UTF-8')
161 data
162
163
164 #Vcorpus(회발성의 corpus를 생성 )를 써서 VectorSource(벡터소스) 를 집합으로 만들어서 넣어준다. 말모음 문치로 하는 작업
165 corpus <- VCorpus(VectorSource(data))
166 corpus
167
168
169 # tm_map(TDM, function) - 전처리를 도와주는데 공백/숫자/특수문자 를 제거하는 함수를 호출하면 처리가 이루어진다.
170
171 # stopwords(불용어) 처리 방법
172 stopwords('en')
173 stopwords2 <- stopwords('en')
174 stopwords2 <- c(stopwords('en'), "and", "not", "but") #더 추가하고 싶을때
175
176 corpus_map <- tm_map(corpus,removeWords,stopwords2) #tm_map (데이터 , 함수(단어를 지우겠다.) , 지울 단어들 (stopword2를 통해 만든 불용어 단어들))
177
178 corpus_map <- tm_map(corpus_map,stripWhitespace) #여러개의 공백을 하나의 공백으로 만들어준다.
179 corpus_map <- tm_map(corpus_map,removeNumbers) #숫자제거
180 corpus_map <- tm_map(corpus_map,removePunctuation) #특수문자 제거
181 corpus_map
182 |
183
184 #빈도 체크
185 TDM <- TermDocumentMatrix(corpus_map)
186
187 findFreqTerms(TDM,2)
188
189 #빈도표
190 matrix <- as.matrix(TDM)
191 rownames(matrix)
192
193 rowSums(matrix)
194
195 wFreq <- sort(rowSums(matrix),decreasing =T)
196
197
198 pal <- brewer.pal(6,"Accent")
199
200 wordcloud(words = names(wFreq),
201           freq = wFreq,
202           min.freq = 1,
203           random.order = F,
204           colors=pal)
205

```

이걸 하나씩 풀어서 설명해보자

```

78 # 텍스트 마이닝 단계
79 # text → corpus(말뭉치) → TDM(term documnet matrix) → TM 분석 → Matrix(DF) → wordcloud
80 # service_data_I_love_mom.txt
81
82 #데이터 읽어오기
83 data <- readLines(file.choose(),encoding = 'UTF-8')
84 data

```

```

> data
[1] "Mommy I love you And I want to say thank you" "I am happy everyday for the love you give to me" "Sometimes I am sad Sometimes I do cry"
[4] "Then I just remember mommy how much you love me" "Mommy you take care of me you help me to be strong" "You help my dreams to come true with everything you do"
[7] "Your love always makes me smile I want to say thank you" "I really, really love you" "Mommy I love you ~~~~~!" "$ I'm

```

#VCorpus - 받은 text를 말뭉치로 만드는 함수

받은 데이터를 말뭉치로 만드는 작업이다. 이렇게 만들어 놔야 전처리를 하기 쉽다.

```

85
86 #Vcorpus(회발성의 corpus를 생성 )를 써서 VectorSource(벡터소스) 를 집합으로 만들어서 넣어준다. 말모음 문치로 하는 작업
87 corpus <- VCorpus(VectorSource(data))
88 corpus
89

```

#TermDocumentMatrix - 행이 단어, 열이 문서 형태로 된 행렬을 만드는 함수

이 말뭉치들을 TermDocumentMatrix를 통해 데이터를 확인한다.

```
> # TermDocumentMatrix : 행이 단어 , 열 문서 형태로 된 행렬
> TDM <- TermDocumentMatrix(corpus)
> TDM
<TermDocumentMatrix (terms: 37, documents: 9)>
Non-/sparse entries: 56/277
Sparsity : 83%
Maximal term length: 10
Weighting : term frequency (tf)
```

9개의 도큐먼트중에 37개의 단어를 를 찾았다.

행열로 봤을때 0 으로 매칭되지 않는 값이 56개가 있다는 뜻이다 . 즉 한 번이라도 노출된 개수가 56개이다.

Maximal term Length 는 가장 긴 단어가 10글자이다.

Sparsity는 희소행렬이다.

#tm_map - 전처리를 도와주는 함수

```
> # tm_map(TDM,function) - 전처리를 도와주는데 공백/숫자/특수문자 를 제거하는 함수를 호출하면 처리가 이루어진다.
> corpus_map <- tm_map(corpus,stripWhitespace) #여러개의 공백을 하나의 공백으로 만들어준다.
> corpus_map <- tm_map(corpus,removeNumbers) #숫자제거
> corpus_map <- tm_map(corpus,removePunctuation) #특수문자 제거
> corpus_map
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 9
> corpus_map[[9]]$content
[1] "Mommy I love you Im"
```

이제 데이터에서 불용어(쓰이지 않는 단어)를 처리한다.

```
176 # tm_map(TDM,function) - 전처리를 도와주는데 공백/숫자/특수문자 를 제거하는 함수를 호출하면 처리가 이루어진다.
177
178 # stopwords(불용어) 처리 방법
179 stopwords('en')
180 stopwords2 <- stopwords('en')
181 stopwords2 <- c(stopwords('en'),"and","not","but") #더 추가하고 싶을때
182
183 corpus_map <- tm_map(corpus,removeWords,stopwords2) #tm_map (데이터 , 함수(단어를 지우겠다.), 지울 단어들 (stopwords2를 통해 만든 불용어 단어들))
184
185 corpus_map <- tm_map(corpus_map,stripWhitespace) #여러개의 공백을 하나의 공백으로 만들어준다.
186 corpus_map <- tm_map(corpus_map,removeNumbers) #숫자제거
187 corpus_map <- tm_map(corpus_map,removePunctuation) #특수문자 제거
188 corpus_map
```

stopwords('en')에는 영어에서 불용어처리 하는 단어들이 들어있다. 그데 이터에 추가적으로 걸리지 않는 and 나 not 등을 추가로 해 stopwords2를 만들어준다.

그후 데이터를 가공하는 작업을 진행한다.

이런 방식으로 한번에 하는 것도 가능하다.

[illegible]

#확인

```
> # stopword 처리 방법
> stopwords('en')
[1] "i" "me" "my" "myself" "we" "our" "ours" "ourselves" "you" "your" "yours" "yourself" "yourselves" "he"
[15] "his" "him" "himself" "she" "her" "hers" "herself" "it" "its" "itself" "they" "them" "their" "theirs"
[29] "themselves" "what" "which" "who" "whom" "this" "that" "these" "those" "am" "is" "are" "was" "were"
[43] "be" "been" "being" "have" "has" "had" "having" "do" "does" "did" "doing" "would" "should" "could"
[57] "ought" "i'm" "you're" "he's" "she's" "it's" "we're" "they're" "i've" "you've" "we've" "they've" "i'd" "you'd"
[71] "he'd" "she'd" "we'd" "they'd" "i'll" "you'll" "he'll" "she'll" "we'll" "they'll" "isn't" "aren't" "wasn't" "weren't"
[85] "hasn't" "haven't" "hadn't" "doesn't" "don't" "didn't" "won't" "wouldn't" "shan't" "shouldn't" "can't" "cannot" "couldn't" "mustn't"
[99] "let's" "that's" "who's" "what's" "here's" "there's" "when's" "where's" "why's" "how's" "p" "a" "an" "the" "and"
[113] "but" "if" "or" "because" "as" "until" "while" "of" "at" "by" "for" "with" "about" "against"
[127] "between" "into" "through" "during" "before" "after" "above" "below" "to" "from" "up" "down" "in" "out"
[141] "on" "off" "over" "under" "again" "further" "then" "once" "here" "there" "when" "where" "why" "how"
[155] "all" "any" "both" "each" "few" "more" "most" "other" "some" "such" "no" "nor" "not" "only"
[169] "own" "same" "so" "than" "too" "very"

> stopwords2 <- stopwords('en')
> stopwords2 <- c(stopwords('en'), "and", "not", "but") # 더 추가하고 싶을때
> stopwords2
[1] "i" "me" "my" "myself" "we" "our" "ours" "ourselves" "you" "your" "yours" "yourself" "yourselves" "he"
[15] "his" "him" "himself" "she" "her" "hers" "herself" "it" "its" "itself" "they" "them" "their" "theirs"
[29] "themselves" "what" "which" "who" "whom" "this" "that" "these" "those" "am" "is" "are" "was" "were"
[43] "be" "been" "being" "have" "has" "had" "having" "do" "does" "did" "doing" "would" "should" "could"
[57] "ought" "i'm" "you're" "he's" "she's" "it's" "we're" "they're" "i've" "you've" "we've" "they've" "i'd" "you'd"
[71] "he'd" "she'd" "we'd" "they'd" "i'll" "you'll" "he'll" "she'll" "we'll" "they'll" "isn't" "aren't" "wasn't" "weren't"
[85] "hasn't" "haven't" "hadn't" "doesn't" "don't" "didn't" "won't" "wouldn't" "shan't" "shouldn't" "can't" "cannot" "couldn't" "mustn't"
[99] "let's" "that's" "who's" "what's" "here's" "there's" "when's" "where's" "why's" "how's" "p" "a" "an" "the" "and"
[113] "but" "if" "or" "because" "as" "until" "while" "of" "at" "by" "for" "with" "about" "against"
[127] "between" "into" "through" "during" "before" "after" "above" "below" "to" "from" "up" "down" "in" "out"
[141] "on" "off" "over" "under" "again" "further" "then" "once" "here" "there" "when" "where" "why" "how"
[155] "all" "any" "both" "each" "few" "more" "most" "other" "some" "such" "no" "nor" "not" "only"
[169] "own" "same" "so" "than" "too" "very" "and" "not" "but"
```

```
121 #빈도 체크
122 TDM2 <- TermDocumentMatrix(corpus_map)
123 corpus_freq <- as.matrix(TDM2)
```

TermDocumentMatrix를 통해 해당문서에서 단어 발생유무를 0,1로 나타낸다.

그 후 그 값을 매트릭스로 형변환 시킨다.

왼쪽행에 값에는 단어가 들어있고 위에 Docs는 이 단어가 발생한 위치의 순서값에 0 / 1로 값이 존재하는지를 나타낸다.

```
> corpus_freq
      Docs
Terms 1 2 3 4 5 6 7 8 9
always 0 0 0 0 0 0 1 0 0
and    1 0 0 0 0 0 0 0 0
care   0 0 0 0 1 0 0 0 0
come   0 0 0 0 0 1 0 0 0
cry     0 0 1 0 0 0 0 0 0
dreams 0 0 0 0 0 1 0 0 0
everyday 0 1 0 0 0 0 0 0 0
everything 0 0 0 0 0 1 0 0 0
for     0 1 0 0 0 0 0 0 0
give    0 1 0 0 0 0 0 0 0
happy   0 1 0 0 0 0 0 0 0
help    0 0 0 0 1 1 0 0 0
how     0 0 0 1 0 0 0 0 0
just    0 0 0 1 0 0 0 0 0
love    1 1 0 1 0 0 1 1 1
makes   0 0 0 0 0 0 1 0 0
mommy   1 0 0 1 1 0 0 0 1
```

다시 전체로 보면

```
159 #데이터 읽어오기
160 data <- readLines(file.choose(),encoding = 'UTF-8')
161 data
162
163
164 #Vcorpus(회발성의 corpus를 생성 )를 써서 VectorSource(벡터소스) 를 집합으로 만들어서 넣어준다. 말모임 문치로 하는 작업
165 corpus <- VCorpus(VectorSource(data))
166 corpus
167
168
169 # tm_map(TDM,function) - 전처리를 도와주는데 공백/숫자/특수문자 를 제거하는 함수를 호출하면 처리가 이루어진다.
170
171 # stopwords(불용어) 처리 방법
172 stopwords('en')
173 stopwords2 <- stopwords('en')
174 stopwords2 <- c(stopwords('en'),"and","not","but") #더 추가하고 싶을때
175
176 corpus_map <- tm_map(corpus,removeWords,stopwords2) #tm_map (데이터 , 함수(단어를 지우겠다.) , 지울 단어들 (stopword2를 통해 만든 불용어 단어들))
177
178 corpus_map <- tm_map(corpus_map,stripWhitespace) #여러개의 공백을 하나의 공백으로 만들어준다.
179 corpus_map <- tm_map(corpus_map,removeNumbers) #숫자제거
180 corpus_map <- tm_map(corpus_map,removePunctuation) #특수문자 제거
181 corpus_map
182 |
183
184 #빈도 체크
185 TDM <- TermDocumentMatrix(corpus_map)
186
187 findFreqTerms(TDM,2)
188
189 #빈도표
190 matrix <- as.matrix(TDM)
191 rownames(matrix)
192
193 rowSums(matrix)
194
195 wFreq <- sort(rowSums(matrix),decreasing =T)
196
197
198 pal <- brewer.pal(6,"Accent")
199
200 wordcloud(words = names(wFreq),
201           freq = wFreq,
202           min.freq = 1,
203           random.order = F,
204           colors=pal)
205
```




꾸까꾸

혼자 끄적끄적하는 블로그 입니다.