

[R] R로 하는 비정형 데이터 처리 (facebook 데이터를 통한 긍정/부정 나누기) — 나무
늘보의 개발 블로그

노트북: blog

만든 날짜: 2020-10-06 오후 4:43

URL: <https://continuous-development.tistory.com/51?category=793392>

R

[R] R로 하는 비정형 데이터 처리 (facebook 데이터를 통한 긍정/부정 나누기)

2020. 8. 3. 17:52 수정 삭제 공개

비정형 데이터 처리

일단 기본적인 패키지들을 install 하자

```
305 # 비정형 데이터 처리(텍스트 마이닝)
306 # 단어 빈도를 나타내는 시각화(wordcloud, koNLP, tm)
307
308
309 install.packages(c("hash", "tau", "Sejong",
310                   "RSQLite", "devtools", "bit",
311                   "rex", "lazyeval", "htmlwidgets",
312                   "crosstalk", "promises", "later",
313                   "sessioninfo", "xopen", "bit64",
314                   "blob", "DBI", "memoise", "plogr",
315                   "covr", "DT", "rcmdcheck", "rversions"),
316                   type = "binary")
317
318 # github 버전 설치
319 install.packages("remotes")
320 # 64bit 에서만 동작합니다.
321 remotes::install_github('haven-jeon/KoNLP',
322                          upgrade = "never",
323                          INSTALL_opts=c("--no-multiarch"))
324
325 |
326
327 # 감성분석
328 # service_data_facebook_bigdata.txt
329
330 fbook <- file(file.choose(),encoding="UTF-8")
331 fbook_read <- readLines(fbook)
332 head(fbook_read)
333 str(fbook_read)
334
```

여기서는 페이스북의 데이터를 가져왔다. 데이터의 내용은 아래와 같다.
이런 식으로 각 행에 대해서 문장형 데이터가 들어가 있었다.

```
> fbook_read <- readLines(fbook) # 한 문장씩 읽어서 제거
> head(fbook_read)
[1] "스마트 기기와 SNS 덕분에 과거 어느 때보다 많은 데이터가 흘러 다니고 빠르게 쌓입니다. 다음 그림은 2013년에 인터넷에서 60초 동안 얼마나 많은 일이 벌어지는지를 나타낸 그림이다. Facebook에서는 1초마다 글이 4만 천 건 포스팅되고, 좋아요 클릭이 180만 건 발생합니다. 데이터는 350GB씩 쌓입니다. 이런 데이터를 실시간으로 분석하면 사용자의 패턴을 파악하거나 의사를 결정하는 데 참고하는 등 다양하게 사용할 수 있을 것입니다."
```

[2] "빅데이터를 처리하는 프레임워크로 흔히 Hadoop MapReduce를 사용한다. MapReduce는 페타바이트 이상의 데이터를 여러 노드로 구성된 클라우드 환경에서 병렬 처리하는 기법으로, 함수형 프로그래밍에서 일반적으로 사용되는 Map과 Reduce 방식을 사용해 데이터를 처리한다. MapReduce는 대량 데이터를 분산 처리할 수 있는 좋은 기법이지만, 배치 방식으로 데이터를 처리하기 때문에 실시간으로 데이터를 조회하기 어렵다. 이런 단점을 극복하기 위해 최근 몇 년간 실시간 분산 쿼리나 스트리밍 처리 기법이 많이 연구되었다."

[3] "실시간 분산 쿼리는 클러스터를 구성하는 노드가 각자 쿼리를 처리하게 해(push down) 한 번에 처리할 데이터의 크기는 작게 하면서 이를 병렬 처리해 응답 시간을 실시간 수준으로 높이는 방식이다. Dremel의 논문을 기반으로 한 Cloudera의 Impala와 Apache Tez, 그리고 최근 공개된 Facebook의 Presto가 이 방식에 속한다."

이 문장들을 전처리할 필요가 있었다. 그래서 정규표현식을 통해 전처리를 하였다.

```
336 #2. 전처리(정규표현식의 필요하다)
337 #문장 부호 제거[[:punct:]]하는 정규표현식 활용
338 #특수문자 제거[[:cntrl:]]
339 #숫자 제거 [[0-9]] \d+(숫자) , \w(단어) , \s+(공백) , \n , \t
340 #gsub() 함수를 이용해서 전처리를 한다.
341
342
343 s1 <- gsub('[[[:punct:]]]', '', fbook_read) #문장부호 제거
344 s1
345
346 s2 <- gsub('[[[:cntrl:]]]', '', s1) # 특수문자 제거
347 s2
348
349 s3 <- gsub('\d+', '', s2) # 숫자제거
350 s3
351
352 s4 <- tolower(s3) #소문자로 변환
353 s4[1]
354
```

gsub을 통해서 정규표현식에 해당하는 데이터를 ""로 변환해주었다.
처음에는 문장부호를 제거하고 그다음에는 특수문자, 숫자제거 이렇게 하였고 마지막에는 모든 대문자를 소문자로 바꿔주는 tolower를 사용하였다.

```
355
356 wordList <- str_split(s4, "\s+") #공백으로 분류
357 wordVec <- unlist(wordList) # vector로 만든다.
```

그다음은 str_split를 사용하여 공백을 통한 단어 분리를 하였다. 그다음 데이터 프레임 형식으로 된 데이터를

[[74]]									
[1]	"빅데이터란"	"엄청나게"	"데이터의"	"양이"	"방대한"	"종래의"	"방법으로는"	"수집"	
[12]	"어려운"	"것을"	"말한다"	"이차적으로는"	"그런"	"단"	"데이터를"	"어려"	
[23]	"정보로"	"만들어내는"	"과정까지를"	"포함한다"	"난"	"강박기부터"	"우리나라에서도"	"무한경쟁의"	
[34]	"새로운"	"플랫폼으로"	"빅데이터란"	"말이"	"대응하기"	"사자했다"	"삼성경제연구소는"	"난"	
[45]	"빅데이터를"	"지속하며"	"이것이"	"미래의"	"성장"	"동력이"	"될"	"거라고"	
[56]									
[[75]]									
[1]	"그런데"	"사실상"	"빅데이터는"	"오래전부터"	"우리"	"실제"	"이미"	"들어와"	"있는"
[14]	"책의"	"저자는"	"말한다"	"또한"	"빅데이터"	"시대에"	"해독능력을"	"위한"	"등계적"
[27]	"데이터를"	"모아"	"분석해"	"가장"	"올바르고"	"빠른"	"답을"	"알려주는"	"실용적인"
[40]	"근거가"	"되기에"	"원래"	"비즈니스맨이"	"지녀야"	"알"	"최강의"	"무기라는"	"것이다"
[[76]]									
[1]	"이"	"책은"	"일본에서"	"통계"	"관련"	"서적으로는"	"이례적으로"	"출간"	"개괄"
[14]	"이례적인"	"현상을"	"물리일키며"	"상반기"	"경제경영"	"분야"	"베스트셀러"	"위에"	"올랐다"
[27]	"적도로"	"통계의"	"역할을"	"새롭게"	"인식한"	"이"	"책은"	"통계학을"	"공부하려는"
[40]	"지금"	"이"	"순간"	"당신의"	"업무에"	"기업에"	"숙한"	"공동체에"	"업무"
[53]	"계획할"	"수"	"있게"	"하는"	"최고의"	"활용서이다"			

unlist를 사용하여 vector 형식으로 바꿔주었다.

[793]	"읽어서는"	"데이터를"	"체계적으로"	"장리하고"	"수집보관할"	"읽"	"필요한"	"데이터를"	"작성하기"
[802]	"접어서"	"사용할"	"수"	"있어야"	"합니다"	"그래서"	"데이터베이스와"	"데이터를"	"관계하고"
[811]	"관리할"	"사용하는"	"사용하는"	"출처에"	"필요합니다"	"데이터를"	"상계적"	"필요"	"사용하기"
[820]	"필요"	"도구도"	"다양하지만"	"지나"	"exerd와"	"sql"	"developer를"	"사용합니다"	"데이터베이스"
[829]	"설계"	"나에게"	"말기"	"exerd"	"exerd는"	"도마도"	"시스템즈에서"	"개발관"	"이클립스"
[838]	"기반의"	"지능을"	"라"	"도구입니다"	"serve"	"무로로"	"이동할"	"수"	"있습니다"
[847]	"관계"	"oracle"	"microsoft"	"sql"	"serve"	"난"	"db"	"mysql을"	"대상으로"
[856]	"리버스프워드"	"엔지니어링과"	"물리적"	"특성"	"권장"	"지향하고"	"있습니다"	"초보지도"	"데이터베이스"
[865]	"설계할"	"직관적이고"	"있고"	"빠르게"	"읽"	"수"	"있습니다"	"--"	"관리"
[874]	"여기"	"하나라"	"중문"	"oracle"	"sql"	"developer"	"sql"	"access에"	"oracle에서"
[883]	"제공하는"	"sql"	"개발도구입니다"	"여러가지"	"데이터베이스와"	"는"	"access에"	"개발"	"데이터베이스"
[892]	"검색할"	"지원합니다"	"데이터베이스에서"	"어라기저"	"권리"	"통제"	"자유자제로"	"데이터를"	"추출하고"
[901]	"가공할"	"수"	"있도록"	"도움을줍니다"	"이유하여"	"대용량"	"실시간"	"데이터와"	"배치"
[910]	"데이터를"	"다양한"	"분석"	"도구들"	"고정해"	"감"	"빠르게"	"분석할"	"수"
[919]	"있는"	"memory"	"computing"	"기반의"	"같이"	"수"	"있을"	"방법적으로"	"빅데이터"
[928]	"시스템을"	"새로운"	"통관"	"대량의"	"데이터"	"수입"	"작업"	"그리고"	"수집된"
[937]	"다양한"	"문제에"	"관련한"	"데이터까지"	"분류하는"	"작업"	"분류한"	"데이터를"	"기반으로"
[946]	"최종"	"결과들을"	"만드는"	"작업으로"	"구성할"	"수"	"있을"	"것이다"	"물론"
[973]	"데이터를"	"기반으로"	"라"	"지"	"데이터를"	"만들"	"수도"	"있고"	"지"
[982]	"데이터를"	"관리"	"관"	"통제"	"데이터를"	"관리"	"분석하여"	"데이터의"	"분석"

이 다음에는 긍정 단어와 부정 단어의 데이터를 통해 내가 가지고 있는 데이터와 매칭 하는 작업을 하였다.

아래 함수를 생성했다. 이 함수를 통해 긍정 / 부정 / 중립을 나눴다.

```

447 library(stringr)
448 library(plyr)
449 ?lapply
450
451 #여기서 생략된 1. list일때 for문처럼 계속 도는것 2.sum(true) 1의 값을 나타내고 3.lapply에서 function을 쓸경우에 매개변수를 밖에서 한번 정의해야 한다.
452
453 # 이 함수를 정의하세요
454 results <- function(words , positive , negative) { #여기서 예초에 리스트로 돌아가기 때문에 for 문처럼 계속 돈다.
455
456   scores = lapply(words, function(words, positive, negative) {
457     pMatch = match(words, positive)
458     nMatch = match(words, negative)
459
460     pMatch = !is.na(pMatch) # true false는 1과 0을 나타낸다. 그래서 이 true를 sum으로 하면 1을 나타내고 이렇게 계산식을 나타낸다.
461     nMatch = !is.na(nMatch)
462
463     score = sum(pMatch) - sum(nMatch) #값이 양수이면 긍정의 단어로 -일경우 부정의 단어로 0 일경우 중립단어로서 나타낸다.
464     return(score)
465   }, positive, negative) # 이거 도구는 문법이다. 안에 function(words, positive, negative) 에 쓴거는 밖에서 이원식으로 적어준다. 결과적으로 return하는 score 하
466
467   scores.df = data.frame(score=scores , text=words)
468   return(scores.df)
469 }
470
471
472 resultTbl <- results(wordVec, pDic, nDic)
473 head(resultTbl)
474
475

```

```

477 resultTbl <- resultS(wordVec, pDic, nDic)
478 str(resultTbl)
479 head(resultTbl)
480
481 resultTbl$text
482 resultTbl$score
483 resultTbl$remark[resultTbl$score ≥ 1] <- "긍정"
484 resultTbl$remark[resultTbl$score = 0] <- "중립"
485
481:1 (Top Level)

```

Console Terminal × Jobs ×

~/ ➔

```

489 0 중분합니다
490 0 스프레드시트
491 0 프로그램은
492 0 셀에
493 0 데이터를
494 0 입력하는
495 0 데
496 0 엑셀
497 0 에서는
498 0 약
499 0 만개 행
500 0 x
[ reached 'max' / getOption("max.print") -- omitted 1988 rows ]
> head(resultTbl)
  score  text
1     0 스마트
2     0 기기와
3     0 sns
4     0 덕분에
5     0 과거
6     0 어느

```

이런식으로 각 단어에 대한 긍정 부정을 볼 수 있다.

```

472
473 # 긍정부정에 따라서 파이차트 만들기
474 resultTbl <- resultS(wordVec, pDic, nDic)
475 str(resultTbl)
476 head(resultTbl)
477
478 resultTbl$text
479 resultTbl$score
480 resultTbl$remark[resultTbl$score ≥ 1] <- "긍정"
481 resultTbl$remark[resultTbl$score = 0] <- "중립"
482 resultTbl$remark[resultTbl$score < 0] <- "부정"
483
484 resultTbl$remark
485
486 table(resultTbl$remark)
487
488 pieResult <- table(resultTbl$remark)
489 pieResult <- table(resultTbl$remark)
490
491 ?pie
492 pie(pieResult,
493     labels=names(pieResult),
494     col = c('yellow', 'green', 'blue'))

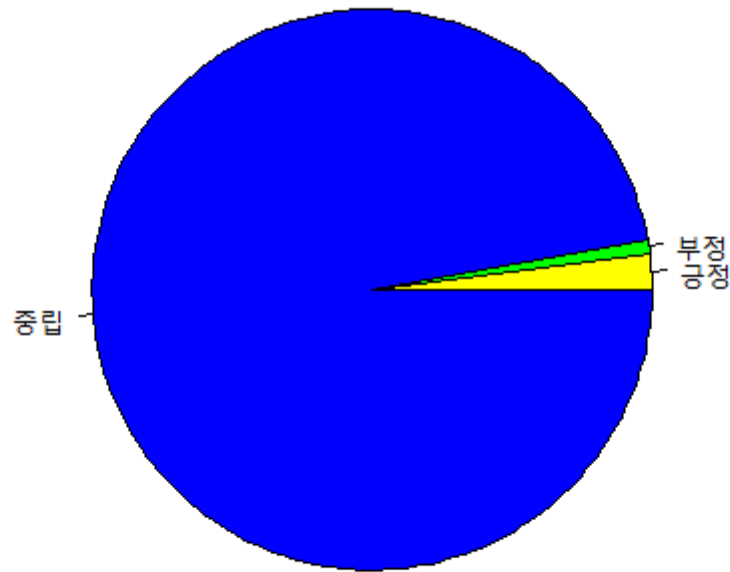
```

```
> table(resultTbl$remark)
```

긍정	부정	중립
51	20	2417

이렇게 긍정 부정중립의 개수를 셀 수 있다.

이걸 파이차트로 나타내면 아래와 같다.



'R' 카테고리의 다른 글

[R] R을 활용한 상관분석과 회귀분석 - 1

[R] R을 통한 텍스트마이닝에서 워드클라우드 까지

[R] R로 하는 비정형 데이터 처리 (facebook 데이터를 통한 긍정/부정 나누기)

[R] R에서 Database 사용하기 / DB 기본적인 구문 사용하기

[R] 예제를 통한 데이터 전처리 작업

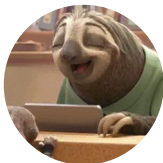
[R] R을 통해 시계열 그래프 만들기 (자료 분석을 위한 시각화와 실습 예제)

R 비정형 데이터 처리

R로 하는 비정형 데이터 처리

비정형 데이터

비정형 데이터 처리



꾸까꾸

혼자 끄적끄적하는 블로그 입니다.