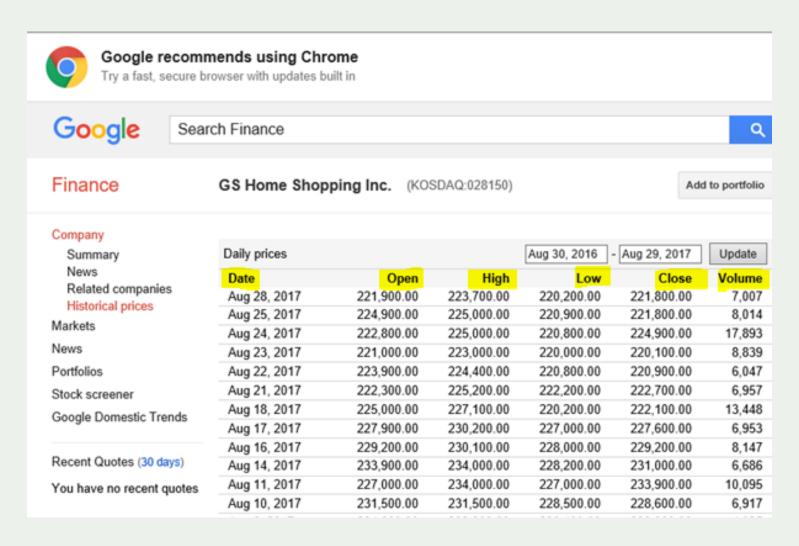
뉴스 텍스트 마이닝

Summary

- 1. Google finance gs 홈쇼핑 주식데이터 크롤링
- 2. Date별 주식이 오르면 up 내리면 down표시
- 3. Daum뉴스에서 gs 홈쇼핑 뉴스데이터 크롤링
- 4. 주식, 뉴스데이터 table merge
- 5. 테이블 결측치 처리
- 6. 웹에서 긍/부정 사전 다운로드
- 7. 크롤링한 뉴스데이터와 비교하여 긍정 사전에 해당하는 단어가 많으면 긍정기사 부정단어에 해당하는 단어가 많으면 부정적 기사로 분류
- 8. 주식 up down과 뉴스의 긍/부정 비교

구글 finance 자료 크롤

-2015/08 ~ 20017/08 3년치 주식데이터 크롤링



R코드- 구글finance data 크롤링

#google finance data 크롤링

```
#https://www.google.com/finance/historical?cid=821222166553443&startdate
= Apr+7 %2C+ 2015 &enddate= May+3 %2C+ 2017 &num= 200
#https://www.google.com/finance/historical?cid=821222166553443&startdate
= Jan+1 %2C+ 2017 &enddate= May+17 %2C+ 2017 &num= 200
startdate < - "2015/08/10"
enddate<-"2017/08/12"
month<-
c("Jan","Feb","Mar","Apr","May","Jun","Jul","Aug","Sep","Oct","Nov","Dec")
startyear < - as.numeric(str_sub(startdate, 1, 4))
startmonth < -month[as.numeric(str_sub(startdate, 6, 7))]
startday < -as.numeric(str_sub(startdate, 9, 10))
endyear < -as.numeric(str_sub(enddate, 1, 4))
endmonth<-month[as.numeric(str_sub(enddate,6,7))]
endday < -as.numeric(str_sub(enddate, 9, 10))
num<-200
start<-0
```

```
url<-"https://www.google.com/finance/historical?cid=821222166553443&startdate="
continue <- TRUE
table < -data.frame()
while(continue)
 url1<-
url%>%paste(.,startmonth,"+",startday,"%2C+",startyear,"&enddate=",endmonth,"+",e
ndday,"%2C+",endyear,"&num=",num,"&start=",start,sep="")
 table_url<-read_html(url1)%>%html_nodes("div#prices")%>%html_nodes(".gf-
table.historical_price")
 plus < -html_table(table_url)[[1]]
 for(i in 2:6) plus[,i] <-as.numeric(gsub(",","",plus[,i]))</pre>
 table < -rbind(table, plus)
 if(nrow(plus) < 200) continue < - FALSE
 start<-start+num
table_month < -factor(str_sub(table$Date,1,3),levels = month,labels=1:12)
table_day<-str_sub(table$Date,5,6)
table_day<-ifelse(str_sub(table_day,2)==",",str_sub(table_day,1,1),table_day)
table_year < -str_trim(str_sub(table$Date,8,12))
```

```
date < -as. Date(paste(table_year, table_month, table_day), format = "%Y%m%d") table_result < -cbind(date, table[, 2:6]) date 20170811 < -c(table_result$Close, 0) date 20170810 < -c(0, table_result$Close)
```

#updown 컬럼 만들기 -

```
updown<-ifelse(date20170810-date20170811>0,"up","down")
updown[1]<-NA
updown<-updown[1:496]
updown<-factor(updown,levels=c("down","up"))
table_result<-cbind(table_result,updown)
table_result
```

구글 finance 주식 자료 크롤링

```
> table_result
                                Low Close Volume updown
          date
                        High
                 Open
    2017-08-11 227000 234000 227000 233900
                                            10095
                                                     < NA >
    2017-08-10 231500 231500 228500 228600
                                              6917
                                                       up
    2017-08-09 231800 232900 230100 230800
                                             4025
                                                     down
    2017-08-08 228200 234500 228200 233900
                                             6076
                                                     down
    2017-08-07 232000 232000 228600 230500
                                             6723
                                                       up
    2017-08-04 229400 237000 228100 231600
                                            15969
                                                     down
7
    2017-08-03 229600 233900 227100 228000
                                            10975
                                                       up
    2017-08-02 233800 233800 228700 228800
                                            10308
                                                     down
    2017-08-01 236100 236100 231500 232700
                                             8863
9
                                                     down
    2017-07-31 232000 236600 225500 235500
10
                                            13809
                                                     down
    2017-07-28 237800 240000 234800 234800
11
                                            12691
                                                       up
12
    2017-07-27 238300 240000 235600 239500
                                            11392
                                                     down
13
    2017-07-26 234900 237800 231300 237700
                                            12830
                                                       up
    2017-07-25 234400 234500 229500 233900
14
                                              8007
                                                       up
    2017-07-24 229000 232600 227300 232100
                                            12644
15
                                                       up
   2017-07-21 228200 231900 226900 227300
                                              9691
16
                                                       up
17
    2017-07-20 229900 230100 226800 226800
                                            10929
                                                       up
    2017-07-19 225000 229800 223800 229800
                                              9571
18
                                                     down
    2017-07-18 225200 225400 223200 225200
                                              9135
19
                                                       up
    2017-07-17 228500 228500 223900 225100
                                            12605
20
                                                       up
    2017-07-14 223800 230000 223700 228500
                                            13138
21
                                                     down
22
    2017-07-13 220000 227000 220000 225900
                                            14706
                                                       up
   2017-07-12 222700 222700 218100 220000
23
                                            14797
                                                       up
    2017-07-11 223500 223900 222500 222700
24
                                            12494
                                                     down
    2017-07-10 225000 226000 222200 223000
                                              7209
                                                     down
```

뉴스 크롤링

-2015/08 ~ 2017/08 3년간 GS홈쇼핑의 뉴스자료 크롤링

А	В	С	D	E	F	G	Н		J	K	L	М	N	0	Р	Q	R
f1	url	date	title	article													
					<mark>은 코렐 비</mark> 긴												
					봇청소기 어												
2	http://new	######	청와대 참	.청와대 참	모진이 보유	구식을 마	각한 것으	로 나타났다	l. 1급 이상	고위공직자	나는 업무 관	<u></u> 년 기업의	주식을 보	유할 수 없	E록 공직자	윤리법이	규정하고
3	http://new	######	탈모닷컴	탈모방지	에 도움을 주	는 스테디	셀러 TS샴푸	푸를 선보이	는 탈모닷컴	섬(대표 장기	영)은 TS심	· 구 1000억	원 판매와	100만 고객	을 돌파했다	구고 7일 밝	혔다.개
4	http://stoc	######	나노스 어	가 5거래	일째 급등하	며 시가총의	성 8위로 올	라섰다.19일	실 오전 9시7	'분 현재 는	전날보다	3650원(25.	26%) 오른	1만8100원	에 거래중이	다. 지난 1	3일 거
5	http://new	######	허창수 "불	"지난 일을	을 잊지 않고	잘 살펴서	일의 지침.	으로 삼는다	ŀ."허창수 (S그룹 회장	t(사진)이 1	9일 서울 역	벽삼동 GSE	[†] 워에서 열	린 3분기 G	S 임원 모임	님에서 [
6	http://new	######	GS '열린 :	GS그룹은	구성원들0	원활하게	소통할 수	있는 '열린	조직문화'	정착에 힘쓰	'고 있다. 일	일과 삶의 조	화를 통해	조직의 활	력과 생산성	은 물론 개	인의 싪
7	http://new	######	집밥같은	40~50대	주부들의 기	·정간편식(ŀ	HMR) 구매	가 빠르게	늘고 있다.	제품 선택의	부이 넓어	진 데다 맛	도 좋아지	면서 1인가	구나 맞벌0	부부뿐 아	니라 직
8	http://stoc	######	코스피 外	,코스피지:	수가 사상 차	음으로 7기	월 연속 성	남승했다.30	일 코스피지	수는 전날	대비 0.169	6 하락한 23	391.79를 기	록했다. 이	날 하락했지	1만 이달 코	크스피지
9	http://stoc	######	코스피 23	(코스피지:	수가 2380선	으로 밀려니	- 하락 흐름	름을 이어기	고 있다.3일	오후 1시	22분 현재 :	코스피지수	는 전날보다	+ 4.36포인.	트(0.18%) ^L	내린 2387.4	3에 거i
10	http://stoc	######	코스피 24	코스피가	개인의 매수	:세에 힘입	어 이틀째 2	2410선에서	상승세를	지속하고 🤉	있다. 와 하여	기닉스는 나	란히 오르.	고 있다.14일	일 오전 9시	13분 현재 :	코스피기
11	http://new	######	GS 허창수	허창수 G	S 회장은 "고	1객 니즈의	변화를 빠.	르게 파악히	나고 유연하	게 대응해 기	지속적으로	고객에게	새로운 가치	l를 제공해	야 하며 내!	부적으로는	원가 ?
12	http://stoc	######	브이티 코	:지엠피의	자회사 브0	티 코스메	틱은 지난 '	1일 GS홈쇼	:핑에서 방송	송한 'VT 블	루 콜라겐	팩트'가 연극	수 완판 행진	일을 이어갔	다고 3일 밝	i혔다VT 블	루콜리
13	http://stoc	######	코스피 9일	코스피지:	수가 9거래?	일 만에 하락	t했다.이로	써 코스피의	기'최장 연속	÷' 사상 최고	그치 기록(8	거래일) 경·	신도 다음의	으로 미뤄졌	다. 코스피의	의 기존 연속	· 최고기
14	http://stoc	######	코스피 장	코스피지:	수가 개장 직	후 사상 최	고치 기록	을 갈아치운	위 숨을 그	그르고 있다	. 오후 들0	서서 지수는	외국인과	기관의 공병	, 탓에 보힌	권을 맴돌.	고 있는
15	http://new	######	화제의 '자	요즘처럼	더운 여름 팀	발뒤꿈치는	여름 샌들	을 신기 전	필수적으로	. 챙겨야 할	부분이다.	허옇게 일(거난 각질은	산들이나	슬리퍼를 신	신으면 더욱	도드리
16	http://new	######	GS홈쇼핑	GS홈쇼핑	이 코렐·비전	선 파이렉스	등으로 유	명한 키친워	어 제조사	월드키친0	투자를 진	!행했다.GS	홈쇼핑은 1	2일 미국계	사모펀드	인 코넬캐피	털과 월

뉴스 크롤링

-2015/08 ~ 2017/08 3년간 주식데이터와 뉴스데이터를 date 기준으로 table merge

	date [‡]	Open [‡]	High [‡]	Low [‡]	Close [‡]	Volume	tom ‡	f1 [‡]	url	title
306	2016-05-10	190800	198900	190500	194000	23175	down	326	http://news.hankyung.com/article/2016051096251	70번의 실패서 배웠다
307	2016-05-11	194100	196900	190300	192900	15056	up	318	http://news.hankyung.com/article/201605112665i	편의점 실적 악화에도
310	2016-05-12	192000	197000	192000	197000	17336	down	320	http://news.hankyung.com/article/2016051263601	해외서 파는 상품에도
309	2016-05-12	192000	197000	192000	197000	17336	down	317	http://news.hankyung.com/article/2016051263681	"중기 사장이 비행기
317	2016-05-18	189800	191900	188300	191000	7443	down	344	http://news.hankyung.com/article/201605181284i	꿰매지 않은 옷 인기
320	2016-05-19	190500	190900	187100	190000	7625	up	346	http://news.hankyung.com/article/201605194619g	세원ITC "사이먼스캇
324	2016-05-25	189100	191000	187000	187100	10538	down	343	http://news.hankyung.com/article/201605250857a	무더위/미세먼지로 0
327	2016-05-30	NA	NA	NA	NA	NA	NA	345	http://news.hankyung.com/article/201605290285i	'저가' 이미지 벗은 홈
331	2016-06-01	185300	186500	181400	185500	6728	up	341	http://news.hankyung.com/article/201606018268g	'정용진표 집밥' PB ፲
330	2016-06-01	185300	186500	181400	185500	6728	up	340	http://news.hankyung.com/article/201606019567i	이마트 롯데에서 PB :
334	2016-06-07	NA	NA	NA	NA	NA	NA	328	http://news.hankyung.com/article/2016060603191	편의점에 IoT 기술 접
336	2016-06-07	NA	NA	NA	NA	NA	NA	357	http://news.hankyung.com/article/2016060608851	일본 수출길 막힌 중기

테이블 결측 처리

-주가데이터 없는 경우 전날 주가로 날짜를 바꿔 해당날짜 이후로 주가를 예측할 수 있도록 함.

	date ‡	Open [‡]	High $^{\updownarrow}$	Low [‡]	Close [‡]	Volumê	tom $^{\updownarrow}$	f1 [‡]	url ÷	title
1	2015-08-10	199100	199100	193100	194500	8579	down	NA	NA	NA
2	2015-08-11	194000	197000	191000	191000	9226	down	NA	NA	NA
3	2015-08-12	191100	194000	188600	189800	18388	8 down		NA	NA
4	2015-08-13	191000	191000	187200	188000	17768	up	NA	NA	NA
5	2015-08-14	NA	NA	NA	NA	NA	NA	555	http://news.hankyung.com/article/201508122585g	홈쇼핑 이젠 홈쇼핑 아닌 '모바
6	2015-08-17	186200	190800	186000	190000	14667	down	554	http://news.hankyung.com/article/2015081609061	사이판 크루즈 탈까'아시아의

- 4	V *												
	date [‡]	Open [‡]	High $^{\updownarrow}$	Low [‡]	Close [‡]	Volumê	tom [‡]	f1 [‡]	url	title			
1	2015-08-10	199100	199100	193100	194500	8579	down	NA	NA	NA			
2	2015-08-11	194000	197000	191000	191000	9226	down	NA	NA	NA			
3	2015-08-12	191100	194000	188600	189800	18388	down	NA	NA	NA			
4	2015-08-13	191000	191000	187200	188000	17768	up	NA	NA	NA			
5	2015-08-13	NA	NA	NA	NA	NA	NA	555	http://news.hankyung.com/article/201508122585g	홈쇼핑 이젠 홈쇼핑 C			
6	2015-08-17	186200	190800	186000	190000	14667	down	554	http://news.hankyung.com/article/2015081609061	사이판 크루즈 탈까			

뉴스 크롤링

-isnews변수생성 : 해당날짜에 뉴스가 있으면 1 없으면 0 으로 표시

	date [‡]	Open [‡]	High ‡	Low [‡]	Close ‡	Volume [‡]	tom [‡]	isnews [‡]
1	2015-08-10	199100	199100	193100	194500	8579	down	0
2	2015-08-11	194000	197000	191000	191000	9226	down	0
3	2015-08-12	191100	194000	188600	189800	18388	down	0
4	2015-08-13	191000	191000	187200	188000	17768	up	1
5	2015-08-13	191000	191000	187200	188000	17768	up	0
6	2015-08-17	186200	190800	186000	190000	14667	down	1
7	2015-08-18	190100	190700	187600	187600	7938	up	1
8	2015-08-19	187100	206300	186000	206200	36964	down	1
9	2015-08-19	187100	206300	186000	206200	36964	down	1
10	2015-08-19	187100	206300	186000	206200	36964	down	1
11	2015-08-19	187100	206300	186000	206200	36964	down	1
12	2015-08-19	187100	206300	186000	206200	36964	down	1

뉴스기사 긍부정 판단

<한글 긍부정 사전 출처>

-https://github.com/The-ECG/BigData1_1.3.3_Text-Mining/blob/master/dictionary.zip

-Github에서 만들어진 Text-Mining에 필요한 긍 부정사전을 google translation한 자료



뉴스기사 긍부정 판단

- -isnew==1 인것만 가져와서 긍부정 사전과 비교
- -Github에서 가져온 긍부정 사전을 이용하여 뉴스 article에 나타난 긍정단어와 부정단어의 개수 sum
- -긍정단어가 많으면 positive article, 부정단어가 많으면 negative article 로 판단-주식데이터의 up, down 과 뉴스데이터의 positive, negative 비교

>	oos_neg_news_	_num									
	date	0pen	High	Low	Close	Volume	updown	isnews	pos_num	neg_num	pos_ar
1	2015-08-13	191000	191000	187200	188000	17768	up	1	14	8	
2	2015-08-17	186200	190800	186000	190000	14667	down	1	22	14	
3	2015-08-18	190100	190700	187600	187600	7938	up	1	17	14	
4	2015-08-19	187100	206300	186000	206200	36964	down	1	9	11	
5	2015-08-19	187100	206300	186000	206200	36964	down	1	13	14	
6	2015-08-19	187100	206300	186000	206200	36964	down	1	5	4	
7	2015-08-19	187100	206300	186000	206200	36964	down	1	9	16	
8	2015-08-19	187100	206300	186000	206200	36964	down	1	12	14	
9	2015-08-20	206100	206100	197000	197300	13321	down	1	7	8	
10	2015-08-21	188000	198200	187700	197200	12127	down	1	12	20	
11	2015-08-21	188000	198200	187700	197200	12127	down	1	10	14	
12	2015-08-21	188000	198200	187700	197200	12127	down	1	8	11	
13	2015-08-24	192000	195800	189200	191400	14519	down	1	8	16	
14	2015-08-24	192000	195800	189200	191400	14519	down	1	16	7	
15	2015-08-24	192000	195800	189200	191400	14519	down	1	14	11	
16	201E 00 21	102000	100000	100200	101/00	1/510	مسمه	1		6	10