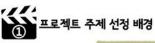






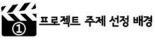
1. 프로젝트 주제 선정 배경



평소 '영화보기'를 좋아하는 A씨. 어느 날, 문득 영화 〈곡성〉이 생각이 났습니다.









(출처 : 네이버 영화(http://movie.naver.com/movie/bi/mi/basic.nhn?code=121051) )

같은 '스토리'와 '반전'에 대한 이야기라도 다른 시각이 존재하고,





(출처 : 네이버 영화(http://movie.naver.com/movie/bi/mi/basic.nhn?code=121051) )

그래서 A씨는 〈곡성〉의 평점을 검색해보았습니다.

- 관람객 평점: 8.23 / 네티즌 평점: 7.60 -

이 점수를 보고 A씨는 '믿고 보기에는 애매'하고 '안 보기는 아쉬운'생각이 들었습니다..

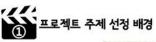
# 프로젝트 주제 선정 배경





(출처 : 네이버 영화(http://movie.naver.com/movie/bi/mi/basic.nhn?code=121051) )

거기다가 댓글은 극과 극으로 나뉘기도 하고,

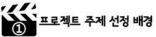


### A씨는 이렇게 평점과 리뷰를 보고 나니 더 혼란스러워졌습니다.





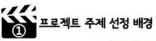






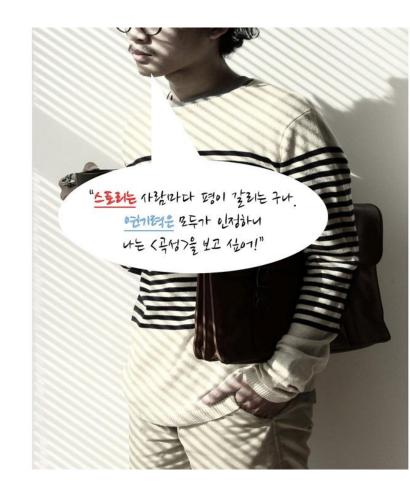
(출처 : 네이버 영화(http://movie.naver.com/movie/bi/mi/basic.nhn?code=121051) )

각 평가 요소별로 다르게 평가하는 글들도 있었습니다.



### 예를 들어.. 이런 결정을 도울 수 있게!

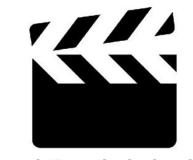




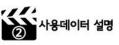


이름하여 맞춤영화리뷰분석기!

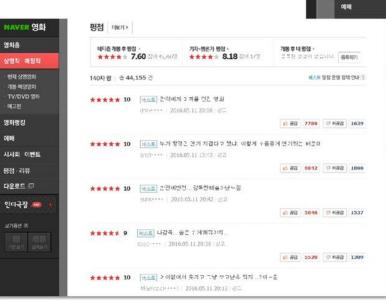




2. 사용 데이터 설명

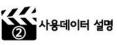


NAVER 영화 〈곡성〉페이지 "140평"



영화랭킹





NAVER 영화

〈곡성〉 페이지 "140평"

별점과 평가글이 함께 작성되는 형태

★★★★★ 10 TV로 결세해서 봤는데 새밌음.

친원덕(xell\*\*\*\*) | 2017.07.01 20:19 |신고

짧은 형태의 글과 길고 구체적인 글이 공<mark>존</mark>한다.

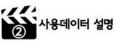
★★★★★ 10 한국의 영화 중 재턴이워스에 가장 근심한 영화이사, 비단 한국 뿐만이 아니라 식이도 40년 온 앞선 영화라고 단연할 수 있다. 모든 대사 하나하나가 의미를 품고 있으면서 시청자의수 준에 따라 미시적으로 보이기도 거시적으로 보이기도 하는 고심리학적영화

아무스(soph\*\*\*\*) | 2017.07.09 05:20 실고

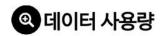
맞춤법에 맞지 않은 글들이 많다.

★★★★★ 1 국성에 대한 명론가들의 명점은 한마디로 여이가 없다. 가가 귤 품어먹는 평점, 아무리 심 오하고 난해한 영화라도 보다보면 일핏 설핏 삼오하고 어렵다는 소공이 가지만 이건 왼전히 끌물이 잡탕 영화 아무 뜻도 없고 그냥 다른 영화 기드라마 짜집?

난해한 영화라도 보다보면 <u>일핏 설핏</u> 심오하고 어렵다는 <u>스긍이.가지만</u>



NAVER 영화 〈곡성〉페이지 "140평"

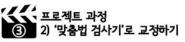


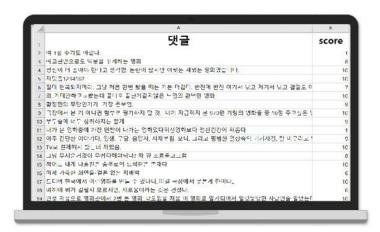
(별점 + 평가글)을 1건으로 했을 때, 약 1만 건.

h	A	8			
1	댓글				
2	야 1점 주기도 아깝다.	1			
3	에고면만으로도 악동을 꾸게하는 영화	U			
4	평점이 다 높아야 한다고 성각함. 논란이 많지만 이쩟든 재밌는 영화였습니다.	10			
5	재밋음123456/	10			
6	절대 현촉되지미리. 그냥 처음 한번 봤을 때는 기분 더럽다. 반전에 반전 여기서 낚고 저기서 낚고 결말도 이	7			
1	와기대안하고ㅠ봤는데 끝나도 끝난거같지않은 느낌의 안백한 영화	10			
8	황성민의 부당연기가 기상 돌보임.	9			
g	극장에서 본 거 아니면 함부로 평가하지 말 것. 내가 지금까지 본 970편 가량의 영화들 중 10집 주고싶은 ﴿	10			
10	부누술에 너무 심취하시는 말게	9			
1	내가 본 영화중에 가장 멘탈이 나가는 영화였다귀신영화보다 정신건강에 허롭다	1			
12	아주 간단한 이야기다. 영생, 무당, 음양사, 시제부림, 오니. 그리고 평범한 일상속의 괴기사건. 잘 버무리고	9			
13	TV로 결제해서 봤는데 재밌음.	10			
14	그냥 무서운거없이 무섭다해야되나? 막 강 소름통고그럼	U			
15	적이도 내게 나홍진은 총구로의 보석같은 존재다	10			
16	어제 가득한 화면들 건은 없는 지배약	6			
17	드디어 한국에서 이런영화를 만들 수 있다니이걸 극장에서 못본게 한이다	10			
18	미끼에 뭐가 걸릴지 모르지만, 새로움이라는 것은 건졌다.	6			
9	인생 처음으로 영화관에서 2번 본 영화, 곽도원을 처음 이 영화로 알게되어서 얼렁뚱땅한 사람인술 알았는데	10			
20	그냥 싫어요이유 없이 싫어요	1			

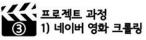


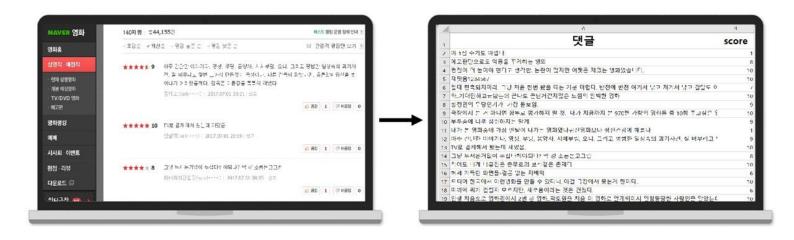
3. 프로젝트 과정





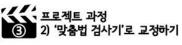
인터넷에서 가져온 글들은 맞춤법 규정에 어긋나는 경우가 많았습니다. 후에 본격적인 텍스트감성분석을 위해서는 맞춤법 검사기를 이용해야 했습니다.





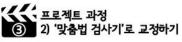
(별점 + 평가글)을 1건으로 했을 때,

약 1만 건 크롤링





NAVER 와 Ddm 의 맞춤법 검사기 중에서 어느 것이 저희가 가지고 있는 텍스트 데이터에 더 알맞는지 테스트 하였습니다.



Before

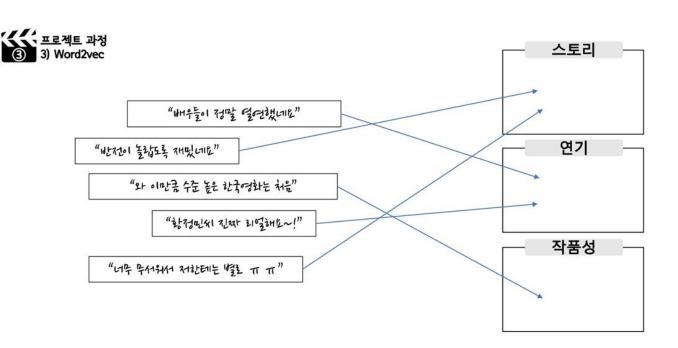
"뒤로가면갈수록 이야기흐름을알게되고 스토리가어떡해흘러가는지 너무뻔했고 개인적으로는 큰흥미는 없엇습니다ㅎㅎ"

After

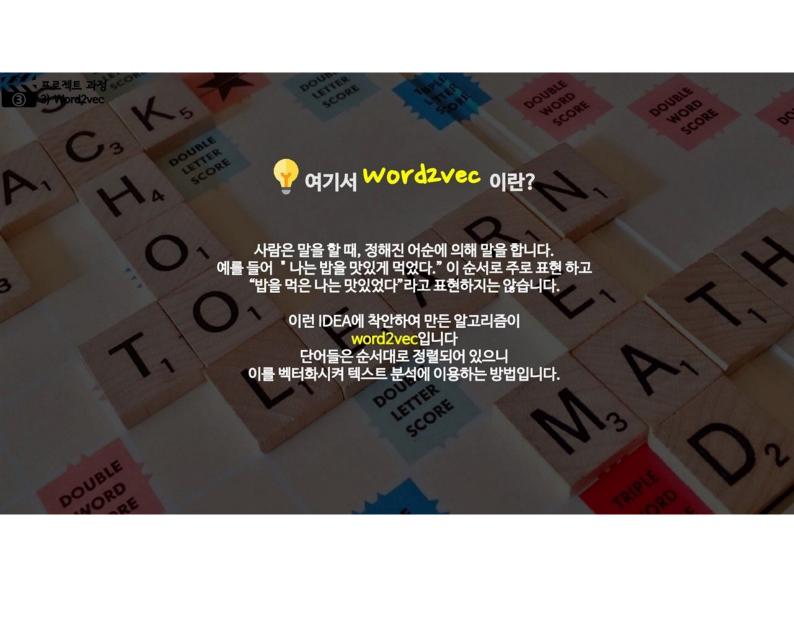
"뒤로 가면 갈수록 이야기 흐름을 알게 되고 스토리가 어떻게 흘러가는지 너무 뻔했고 개인적으로는 큰 흥미는 <mark>없었습니다</mark> ㅎㅎ"

(철자 / 파랑색 :띄어쓰기 / 초록색 :철자+띄어쓰기 / 분홍색 :맞춤법의심)

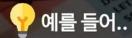
그리고 크롬 드라이버를 다운로드 받아 selenium을 이용하여 10000건의 텍스트가 맞춤법에 맞게 고쳐지도록 하였습니다.



다음으로 다양한 평가요소들을 각각 카테고리별로 분류하기 위해 word2vec 분석을 했습니다.







아래와 같이 word2vec 학습시켜 단어를 벡터화 시킨 matrix가 나왔다고 했을 때

	V1	V2	V3	 V49	V50		
영화	-0.32	0.15	-0.33	0.23	0.34		
생각	-0.12	-0.31	-0.30	0.14	-0.22		
:							
이해	-0.46	0.33	-0.15	0.41	-0.31		

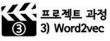
## 코사인 유사도

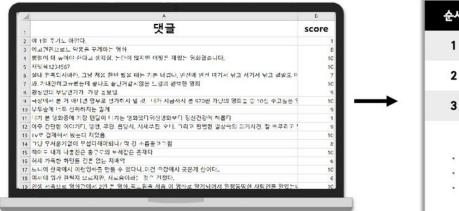
코사인 유사도 
$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum\limits_{i=1}^{n} A_i \times B_i}{\sqrt{\sum\limits_{i=1}^{n} (A_i)^2} \times \sqrt{\sum\limits_{i=1}^{n} (B_i)^2}}$$

- Similarity (영화, 생각) '영화' 벡터와 '생각' 벡터 간 코사인 유사도로 정의 유클리드 거리와는 다르게 벡터의 방향성을 고려한 거리 척도



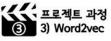
⇒ 문장1은 '스토리' 키워드로 분류!

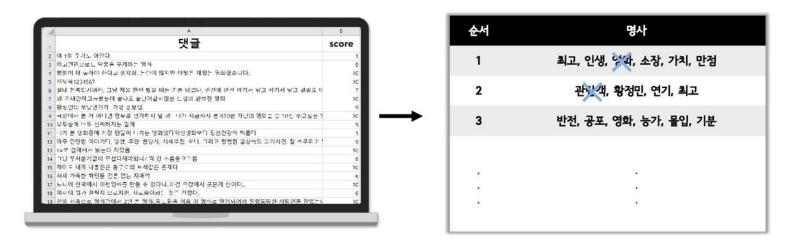




순서	명사
1	최고, 인생, 영화, 소장, 가치, 만점
2	관람객, 황정민, 연기, 최고
3	반전, 공포, 영화, 능가, 몰입, 기분
	8
	#
	<u> </u>

이렇게 word2vec 분석을 하기 위해, 먼저 KONLPY 패키지를 이용해 리뷰 명사 추출하였습니다.

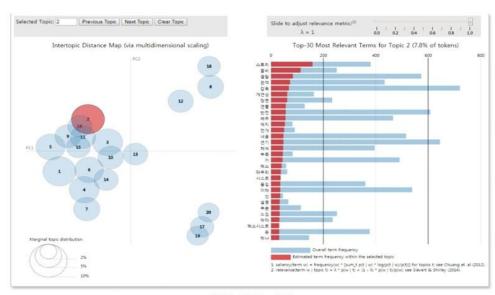




여기서 <u>영화, 관람객</u> 등의 단어는 빈번하게 나타나므로 제거를 해줍니다 -> 다음으로 LDA를 이용해 키워드를 분류하였습니다.

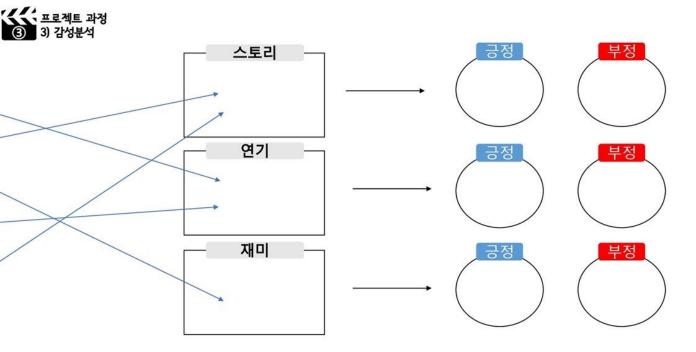






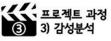
<LDA 분석 예시>

LDA 분석을 통해 word2vec에 활용할 대표성을 띄는 키워드를 살펴보았습니다. 그런 후에 저희는 〈재미〉,〈연기〉,〈이해〉,〈긴장감〉,〈스토리〉를 대표성을 띄는 키워드로 선별했습니다.



이렇게 각 대표적인 키워드에 맞춰 분류한 평<del>들을</del> 이제는 긍정적인 평과 부정적인 평으로 나누기 위해 <mark>감성분석</mark>을 하였습니다.



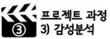


```
[('사랑', 'NNG'), ('하', 'XSV'), ('하', 'XSV'), ('어요', 'FEN'), ('어요', 'FEN'), ('부', 'NNG'), ('박', 'NNG'), ('박', 'NNG'), ('희', 'VV')]
```

<형태소 분석 예시>

여기서 **"투빅스"**처럼 <mark>기존에 없던 단어들은</mark> 알맞은 형태로 나누어지지 않습니다.

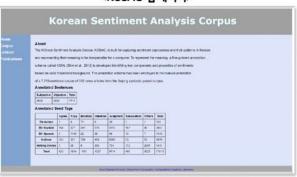
그래서, 이를 보안하고자 감성사전에 신조어가 알맞게 분석될 수 있게 단어를 추가하였습니다.



#### KOSAC(Korean Sentient Analysis Corpus) : 한국어 감정 및 의견 분석 코퍼스

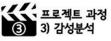
- 단순히 궁&부정 단어의 빈도수를 연산하는 방법에서 벗어나서 포괄적인 감정 표현을 학습시켜 대량의 코퍼스를 구축.
  - \_\_ 서울대 언어학과 연구진(신효필 외 3인)이 한국연구재단의 2년간의 지원을 받아 제작하였다.

#### <KOSAC 홈페이지>



(http://word.snu.ac.kr/kosac/index.php)

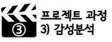
저희는 감성분석을 하기위해 텍스트 데이터에 극성을 달아야 했습니다. 이를 위해, 저희는 KOSAC(한국어 감정 및 의견 분석 코퍼스)를 활용하였습니다.



```
[('사랑', 'NNG'),
('하', 'XSV'),
('하', 'XSV'),
('어요', 'EFN'),
('투', 'NNG'),
('박', 'NNG'),
('박', 'NNG'),
('슬', 'W')]
```

<형태소 분석 예시>

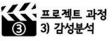
KOSAC 감성사전을 이용하기 위해서는 먼저 <mark>형태소 분석</mark>을 해야합니다. 형태소 분석기는 KOSAC 감성사전의 형식과 가장 유사한 형태로 나눠주는 KONLPY에 내장되어 있는 '꼬꼬마 형태소 분석기'를 활용하였습니다.



1	A	В	C	D	E	F	G	Н	- 1	J
1		Inten-H	expr-l	expr-da	expr-de	expr-ds	expr-wd	inten-L	inten-M	inten-N
2	0	3.796578	9.259459	0.541324	11.45394	11.45561	3.289674	3.578908	25.49008	3.134433
3	1	4.160713	8.909123	0.734143	12.56231	18.10802	5.686405	4.320725	32.06172	5.456844
4	2	2.738262	4.388276	0.011111	5.183975	8.292806	2.123831	2.163845	12.30763	2.790258
5	3	7.735428	8.45102	0.037366	8.395229	17.35531	4.761074	4.15061	22.94254	4.171419
6	4	1.585714	4.773485	0	7.497835	12.33149	2.397186	3.384336	19.67443	2.355519
7	5	0.525499	1.667704	0.006452	0.938978	2.332124	0.054743	0.748272	3.660407	0.065822
8	6	0.501595	1.122678	0.037366	3.631041	4.604915	2.604	1.184917	8.760141	1.553347
9	7	5.942323	5.0947	0.189973	12.62555	19.32692	2.762864	6.251678	24.81878	2.987223
10	9	3.550417	9.567515	0.059938	10.67204	15.41815	4.282362	5.094851	28.15245	3.202287

:

앞서 말한 과정<del>들을</del> 통해, 영화 평<del>들은</del> 총 <mark>29가지 변수를</mark> 가진 데이터가 됩니다.



4	Α	В	C	D	E	F	G	Н	1	
1	Х	review_data	스토리	재미	연기	이해	긴장감	classified	type	
2	similarity	영화를 보는 [	1.392846	1.572395	0.61516	0.751653	1.946108	긴장감	neut	
3	similarity	분위기로 보는	0.825231	1.104023	0.771678	0.694231	0.838665	재미	neg	
4	similarity	뻔할 수 있는	2.481165	1.795907	4.219205	0.826535	1.626136	연기	pos	
5	similarity	무섭지만은 않	1.338093	1.493557	0.857219	1.133663	1.470902	재미	neg	
6	similarity	'전반부는 괜	3.297156	3.326167	4.158264	1.382304	2.665082	연기	neg	
7	similarity	2번 봤는데 2	1.147155	1.970664	0.933674	1.331418	2.188829	긴장감	pos	
8	similarity	배우들의 맛낕	2.736745	2.010845	4.106133	1.585055	2.246647	연기	neg	
9	similarity	시간 가는 줄	3.194446	2.307892	4.08909	1.663725	3.231101	연기	neg	
10	similarity	곡성의 가장 4	1.013364	1.328693	0.933065	0.664448	0.983112	재미	neg	
11	similarity	영화에 호불:	3.008921	3.804916	4.085896	2.91862	2.870648	연기	pos	
12	similarity	내가 생각했던	3.908174	3.224646	4.082342	2.7441	3.480821	연기	neut	
13	similarity	집 가서 자꾸	1.363797	1.840394	1.008006	2.683451	1.405013	이해	neg	

이 데이터를 분석하여 긍정/중립/부정으로 분류하였습니다. 분류기는 여러 분류기를 사용해본 결과 가장 좋았던 Random Forest(약 60%)를 사용하였습니다.



4. 결과 및 제언



**Positive** 

〈재미〉

Negative

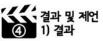
"실거 스웰커, 어적인 소재를 닦는 어뢰를 좋아라는데 친구에서 이렇 작물을 만들 좋아다. 긴 시간이였는데도 지루한 같이 때닷데요"

"아니상아이만 한 주인공들의 관계, 재미보다는 관기를 햇빛되게 하는데 흔 정신이 열대있는 것 같아서 아쉽지만 나를 훌륭한 작물 감습니다."

"2년 발도 본 통에 진이 다 WHAIIII®, 7H관성 파지시는 볼틀 이 정화가 다루는 주제 자체가 카고스합니다. 필신시키는 집은 철근 본 정화 중 최고네요." "무슨 말이지는 알겠는데 그냥 제네가 있다. "

"재미집중은 진짜 기다라고 봤는데 소를 가치고 간자되는 불위기는 작 소식한 것 같으나 별 각종 집중을 "

"곡성 하시기 다이 잇길에 하시기 보고 이하나라고 나'니 저띠있었는 근데 제띠댔지는 땡은 뛰지?"



**Positive** 

(이해)

Negative

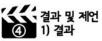
"커말 지극해서 봤어요. 이해 한되는 부별은 다른 불물 줄 쏟게 보구 이해했더했지요. 생각활수록 무서워서 떠칠 잘 설렜네요."

"바라에 바라 전기력도 대단 이느 정도 이해하고 보면 재미있는 영화가 될 듯 해요"

"두 반이나 봤지도 첫 번째 봤을 땐 재밌기는 해도 내용이 잘 이해가 한 가는지라고도 두 번 봐도 가누우 "재밋딧는데 내용이 완전히 이해되지 않았다. "

"볼막성었습니다,,, 물제는 취실부터 끝까지 잘 받아는 이해가 간다는 것"

"도 구고 스트레스바는데 불만한 역할. 두 12은 발아는 되는 역할"



**Positive** 

〈연기〉

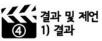
Negative

"무서운 장면 하나 없이 이라이라하게 무섭다. 전혹라 전기, 그리고 쿠너무나 중라 이번 아예내무에게 합도당해버졌다."

"스스타는 벌위기가 작품, 내거기,이따지라 작 맛고 내거우들의 열면이 돋보이다."

"목성 대ዚ 안 1호 아마물이 띴仄던 이항 . 락도원 항정인의 제의 다음에도 또 이렇게 이길," "건기자의 건기가 단면 10점 막겁에 10점 군에 귀시 배기로 불러 들어가는 것 같아 친구이 한되고 군사한 분위기와 배우들의 전기로 돌십되는 영화지 내용은 그냥."

"배수들이 전기를 잘해서 2점이라도 줬다. 이렇게 잔이라고 무서울에 어찌 15세 관站가라는 건지. 항건이의 건체는 무엇이다 천수회의 전체는 뭐라는 건지 . 결약을 관객에게 따루는 전 뭐站, 무서울라 압당함 쩹ೃ함한 넣는 전화."



**Positive** 

〈스토리〉

Negative

"결혼 재밌습니다 영화가 많는 메시지는 진지 참니다 한국적인 해락적인 요소들도 있고 计여든 재밌이다"

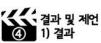
" 바지막 가지 이느 쪽이 문제인지 갈등을 교소시키는 구성이 정말 좋았고 전기자들의 집중적도 대단해서 보는 내내 건강감을 볼 수 집중합니다"

> "언기 물거의 지루창 조차 확이볼 수 때는 전지 수준 또한 대단하다."

"두서 때로 디치스러운 스틸리 도원이 joi 캐스팅 비스인 거 같고 더라운 그래간만에 실막지 이네요"

"이건 뭐 스토니는 날로 쓴 것보다 찾는다 뭐던 대도 아까다"

"소설로 원자이 있을 것 가능을 건도로 스트의 저건내가 좋았습니다 가득의 숨어진 복사를 찾아보는 재미가 꽤 쏠쏠나다군요"



**Positive** 

〈 긴장감 〉

Negative

"실수일 사이에 두 번에나 본 처음 볼 때는 건강간으로 두 번째는 스틸리 100% 이해를 제되가 더 참 "

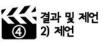
"최고입니다 와 主불구가 결되는지 이해가 한 겁니다 구성엔에서 완벽참니다 추거자 취원, 환지도 진자 10분 월 수 없는 진짜 터지"

> "너무 재밌게 시간 가는 줄 뜨고 봤이요 상이시간이 이렇게 긴 줄 몰랐는데 지루하던 생각이 전혀 안 들었네요"

"소바 코벡라 스킬을 떠나는 해내가 좋았는데, 갈수록 스킬만이 제속되는 건강한데 힘들었답니다 끝나고 나서는 알 수 때는 캠핑함이 났는"

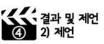
""곱게으면 생활수록 만吹이 나는 생활할 같은 정화십니다"

"전체적인 내용과 진장하는 칼값으나 너무나 많은 복사를 심어놓은 개발에 본데의 스트되에 진출 할 수 없었다"



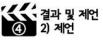
〈맞춤영화리뷰분석기〉 완성! 사실 아쉬운 점이 많다..





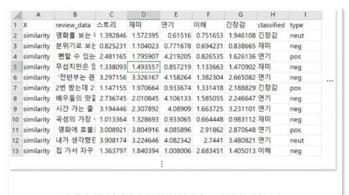
그래도 word2vec을 통해 평가 요소별 카테고리를 나누는 것은 성공적이었고,





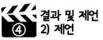
#### 먼저, 감성분석의 정확도가 약 60%밖에 안 되어 실사용은 힘들다는 점..

#### 하지만, 3가지 무작위 분리한다 했을 때 정확도는 33%이며, 우리 모델의 Lift 값이 2이기 때문에 괜찮은 편이다.

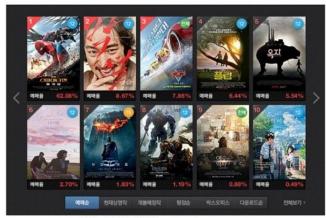


이 데이터를 분석하여 긍정/중립/부정으로 분류하였습니다. 분류기는 여러 분류기를 사용해본 결과 가장 좋았던 Random Forest(약 60%)를 사용하였습니다.





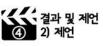
# 그리고 (곡성)에만 적용해본 것이라는 점.



(출처 : 네이버 영화(http://movie.naver.com/))

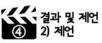
하지만, 이는 다른 영화에 같은 코드를 사용하면, 각 영화 별 분석이 가능해진다.





전체적인 결과물은 정확도와는 별개로 생각보다 괜찮았고,





신조어를 감성사전에 추가하는 등 다양한 시도를 해봤다는 점은 의미 있었다고 자체평가 해본다..ㅎㅎ





