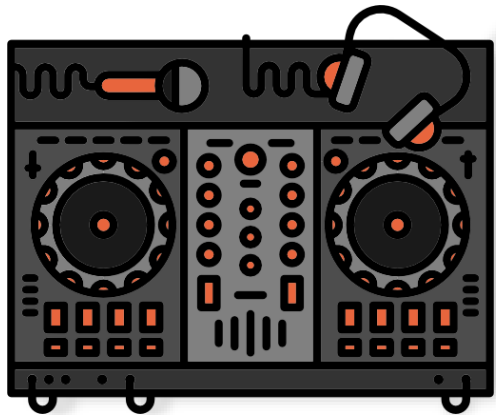


A close-up, low-angle shot of a robotic arm, possibly a DJ mixer or a specialized robotic hand, with numerous white cables attached to its joints and base. The arm is positioned over a computer keyboard, which is visible in the lower-left corner. The background is a solid teal color. The overall lighting is dim, with the white cables providing a strong contrast.

AI DJ

비정형 데이터분석 4조

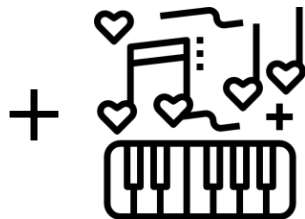
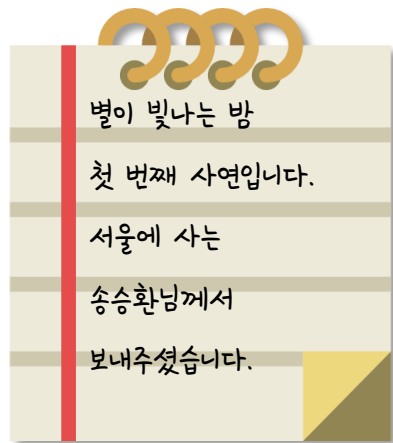
박건빈 변윤선 이수빈



INDEX

01 주제 02 데이터 수집 03 데이터 정제 04 모델링 05 결과

01 주제



대숲 사연(Text)기반
음악 추천 시스템 구현

걱정 고민 많은 20대 청춘들이
사연을 많이 올리는 '대나무숲'

진행과정

데이터 수집 ➤ 데이터 정제 ➤ 모델링 ➤ 추천시스템 구현 ➤ 결론

진행과정

데이터 수집

데이터 정제

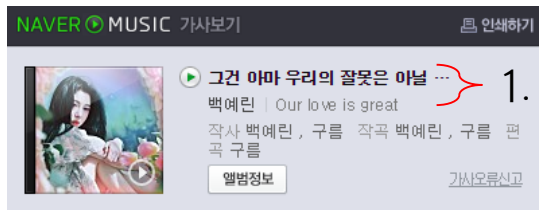
모델링

추천시스템 구현

결론

- 음악 정보데이터
- 서울 20개 대학 대나무숲

-네이버 뮤직 + 벅스 뮤직



사실은 나도 잘 모르겠어
불안한 마음은 어디에서 태어나
우리에게까지 온 건지

나도 모르는 새에 피어나
우리 사이에 큰 상처로 자라도
그건 아마 우리의 잘못된 아날 거야

그러니 우린 손을 잡아야 해
바다에 빠지지 않도록
끊임없이 눈을 맞춰야 해
가끔은 너무 익숙해져 버린
서로를 잃어버리지 않도록

2. 가사

Bugs!

그건 아마 우리의 잘못된 아날 거야



#나들이/소풍 #봄 #지치고 힘들때 #감성적인 #봄비 #화창한 날 #목소리/음색 #지하철/버스 #인디
#취준/수험생 #카페 #학교/퇴근길 #가요 #이색테마 #아티스트 #비/호림 #국내인디 #백예린

3. 테마

- 국내, 주간 TOP 100, 2016년6월-2019년3월 3째주까지, 약 130주 (130주 X 100= 13,000 곡)
- Crawling 항목: 1. 노래제목 2. 노래가사

-네이버 뮤직 + 벅스 뮤직

NAVER MUSIC 가사보기

1. 제목

그건 아마 우리의 잘못된 아닐 ...
백예린 Our love is great
작사 백예린, 구름 작곡 백예린, 구름 편곡 구름

앨범정보 가사오류신고

사실은 나도 잘 모르겠어
불안한 마음은 어디에서 태어나
우리에게까지 온 건지

나도 모르는 새에 피어나
우리 사이에 큰 상처로 자라도
그건 아마 우리의 잘못된 아닐 거야

그러니 우린 손을 잡아야 해
바다에 빠지지 않도록
끊임없이 눈을 맞춰야 해
가끔은 너무 익숙해져 버린
서로를 잃어버리지 않도록

2. 가사

Bugs!

그건 아마 우리의 잘못된 아닐 거야



#나들이/초콜릿 #봄 #지치고 힘들때 #감성적인 #화창한 날 #여름의 추억 #지하철/버스 #인디

#취준/수험생 #카페 #학교/퇴근길 #가요 #이색테마 #아티스트 #비/호림 #국내인디 #백예린

3. 테마

- 국내, 주간 TOP 100, 2016년6월-2019년3월 3째주까지, 약 130주 (130주 X 100= 13,000 곡)
- Crawling 항목: 1. 노래제목 2. 노래가사

데이터 수집

1) 음악 데이터 수집

곡명, 아티스트, 가사

title	artist	lyric	랩, 힙합	휴식, 힐링	카페	비오는 날	발라드	R.B.Soul	사랑, 설렘	여행, 산책
아기공룡 둘리	만화천국	아기공룡 둘리	0	0	0	0	0	0	0	0
로봇 수사관	NA	힘차게 달려	0	0	0	0	0	0	0	0
원숭이 (Pe조권)	말썽꾸러기	말썽꾸러기	0	0.2	0	0	0	0	0	0
날아라 슈퍼	NA	치키치키치	0	0	0	0	0	0	0	0
검정고무신	NA	할아버지	0	0	0	0	0	0	0	0
명탐정 코난	NA	아침에 눈	0	0	0	0	0	0	0	0
우리의 꿈	NA	내 어린 시	0	0	0	0	0	0	0	0.125
아기공룡 둘리	NA	요리보고	0	0	0	0	0	0	0	0
바람의 빛	오연준	사람들만	0	0.428571	0	0	0	0.214286	0.071429	0
요술공주	NA	너와 나의	0	0	0	0	0	0	0	0
달의 요정	NA	미안해 솔	0	0	0	0	0	0	0	0
꾸러기 수호	NA	돌기 명이	0	0	0	0	0	0	0	0
천사소녀	NA	오늘밤에	0	0	0	0	0	0	0	0.111111
고요한밤	동방신기	고요한밤	0	0	0	0	0.375	0	0.125	0
글로리아	동방신기	지극히 높	0	0.060976	0.012195	0	0.04878	0	0.02439	0
The First	동방신기	저 들밖에	0	0	0	0	0.266667	0	0.2	0
크리스마스	러브키즈	크리스마스	0	0	0.2	0	0	0	0.2	0
겨울아이	빅마마	겨울에 태	0	0.047619	0.142857	0	0.238095	0	0.190476	0
소중한 사	길은정	높아만 가	0	0	0	0	0.5	0	0	0

테마

드라이브	OST	재즈	겨울연가	봄의왈츠	여름향기	댄스	가을동화	감성	트로트
0	0.333333	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0.272727	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0.272727	0	0	0	0	0	0	0	0
0	0.25	0	0	0	0	0	0	0	0
0	0.222222	0	0	0	0	0	0	0	0
0	0	0	0.071429	0	0	0	0	0	0
0	0.375	0	0	0	0	0	0	0	0
0	0.4	0	0	0	0	0	0	0	0
0	0.25	0	0	0	0	0	0	0	0
0	0.222222	0	0	0	0	0	0	0	0
0	0	0	0.25	0	0	0.125	0	0	0
0.121951	0	0	0.012195	0	0.207317	0.085366	0	0	0
0	0	0	0.266667	0	0	0.133333	0	0	0
0	0	0	0.6	0	0	0	0	0	0
0	0	0	0.285714	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0.166667
0	0.235294	0	0.117647	0	0	0	0	0	0
0	0	0	0.321429	0	0	0.071429	0	0	0
0	0	0	0.2	0	0	0	0	0	0
0	0.083333	0	0	0	0	0.333333	0	0	0

...

3만곡의 노래 정보

곡명, 아티스트, 가사, 멜론 DJ플레이리스트에 등록 되어있는 테마 데이터를 사용

02

데이터 수집

1) 음악 데이터 수집

곡명, 아티스트, 가사

테마

title	artist	lyric	랩,힙합	휴식,힐링	까페	비오는날	발라드	R.B.Soul	사랑,설렘	여행,산책
아기공룡 하나화천국	아기공룡 하나	아기공룡 하나	0	0	0	0	0	0	0	0
로봇 수사관	로봇 수사관	로봇 수사관	0	0	0	0	0	0	0	0
원숭이 (Pe조권)	원숭이 (Pe조권)	원숭이 (Pe조권)	0	0.2	0	0	0	0	0	0
날아라 슈퍼맨	날아라 슈퍼맨	날아라 슈퍼맨	0	0	0	0	0	0	0	0
검정고무신	검정고무신	검정고무신	0	0	0	0	0	0	0	0
명탐정 코난	명탐정 코난	명탐정 코난	0	0	0	0	0	0	0	0
우리의 꿈	우리의 꿈	우리의 꿈	0	0	0	0	0	0	0	0.125
아기공룡 하나	아기공룡 하나	아기공룡 하나	0	0	0	0	0	0	0	0
바람의 빛	바람의 빛	바람의 빛	0	0.428571	0	0	0	0.214286	0.071429	0
요술공주	요술공주	요술공주	0	0	0	0	0	0	0	0
달의 요정	달의 요정	달의 요정	0	0	0	0	0	0	0	0
꾸러기 수호대	꾸러기 수호대	꾸러기 수호대	0	0	0	0	0	0	0	0
천사소녀	천사소녀	천사소녀	0	0	0	0	0	0	0	0.111111
고요한밤	고요한밤	고요한밤	0	0	0	0	0.375	0	0.125	0
글로리아	글로리아	글로리아	0	0.060976	0.012195	0	0.04878	0	0.02439	0
The First	The First	The First	0	0	0	0	0.266667	0	0.2	0
크리스마스	크리스마스	크리스마스	0	0	0.2	0	0	0	0.2	0
겨울아이	겨울아이	겨울아이	0	0.047619	0.142857	0	0.238095	0	0.190476	0
소중한 사람	소중한 사람	소중한 사람	0	0	0	0	0.5	0	0	0

드라이브	OST	재즈	겨울연가	봄의왈츠	여름향기	댄스	가을동화	감성	트로트
0	0.333333	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0.272727	0	0	0	0	0	0	0	0
0	0.4	0	0	0	0	0	0	0	0
0	0.222222	0	0	0	0	0	0	0	0
0	0	0	0.071429	0	0	0	0	0	0
0	0.375	0	0	0	0	0	0	0	0
0	0.4	0	0	0	0	0	0	0	0
0	0.25	0	0	0	0	0	0	0	0
0	0.222222	0	0	0	0	0	0	0	0
0	0	0	0.25	0	0	0.125	0	0	0
0.121951	0	0	0.012195	0	0.207317	0.085366	0	0	0
0	0	0	0.266667	0	0	0.133333	0	0	0
0	0	0	0.6	0	0	0	0	0	0
0	0	0	0.285714	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0.166667
0	0.235294	0	0.117647	0	0	0	0	0	0
0	0	0	0.321429	0	0	0.071429	0	0	0
0	0	0	0.2	0	0	0	0	0	0
0	0.083333	0	0	0	0	0.333333	0	0	0

과거 멜론에서 크롤링한 약 30,000건의 곡 정보 활용

...

3만곡의 노래 정보

곡명, 아티스트, 가사, 멜론 DJ플레이리스트에 등록 되어있는 테마 데이터를 사용

02

데이터 수집

1) 음악 데이터 수집

곡명, 아티스트, 가사

테마

title	artist	lyric	랩,힙합	휴식,힐링	카페	비오는날	발라드	R.B.Soul	사랑,설편	여행,산책
아기공룡 둘리만화천국	아기공룡 둘리	아기공룡 둘리만화천국	0	0	0	0	0	0	0	0
로봇 수사관	NA	힘차게 달려	0	0	0	0	0	0	0	0
원숭이 (Pe조권)	말썽꾸러기	말썽꾸러기	0	0.2	0	0	0	0	0	0
날아라 슈퍼맨	지키맨	지키맨	0	0	0	0	0	0	0	0
검정고무신	NA	할아빠	0	0	0	0	0	0	0	0
명탐정 코난	아침에 눈	아침에 눈	0	0	0	0	0	0	0	0
우리의 꿈	NA	내 어린 시	0	0	0	0	0	0	0	0.125
아기공룡 둘리	NA	요리보고	0	0	0	0	0	0	0	0
바람의 빛	오연준	사람들만	0	0.428571	0	0	0	0.214286	0.071429	0
요술공주	NA	백과 나의	0	0	0	0	0	0	0	0
달의 요정	NA	백과 나의	0	0	0	0	0	0	0	0
꾸러기 수사관	NA	둘기 명이	0	0	0	0	0	0	0	0
천사소년	NA	오늘밤에	0	0	0	0	0	0	0	0.111111
고요한밤	동방신기	고요한밤	0	0	0	0	0.375	0	0.125	0
글로리아	동방신기	지극히 높	0	0.060976	0.012195	0	0.04878	0	0.02439	0
The First	동방신기	저 들밖에	0	0	0	0	0.266667	0	0.2	0
크리스마스러브키즈	크리스마스	크리스마스	0	0	0.2	0	0	0	0.2	0
겨울아이	빅마마	겨울에 태	0	0.047619	0.142857	0	0.238095	0	0.190476	0
소중한 사람	길은정	높아만 가	0	0	0	0	0.5	0	0	0

드라이브	OST	재즈	겨울연가	봄의왈츠	여름향기	댄스	가을동화	감성	트로트
0	0.333333	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0.272727	0	0	0	0	0	0	0	0
0	0.222222	0	0	0	0	0	0	0	0
0	0	0	0	0.071429	0	0	0	0	0
0	0.375	0	0	0	0	0	0	0	0
0	0.4	0	0	0	0	0	0	0	0
0	0.25	0	0	0	0	0	0	0	0
0.121951	0	0	0.012195	0	0.207317	0.085366	0	0	0
0	0	0	0.266667	0	0	0.133333	0	0	0
0	0	0	0.6	0	0	0	0	0	0
0	0	0	0.285714	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0.166667
0	0.235294	0	0.117647	0	0	0	0	0	0
0	0	0	0.321429	0	0	0.071429	0	0	0
0	0	0	0.2	0	0	0	0	0	0
0	0.083333	0	0	0	0	0.333333	0	0	0

과거 멜론에서 크롤링한 약 30,000건의 곡 정보 활용

→ 중복 곡 제거 후 총 32,000곡 음악 데이터 확보

3만곡의 노래 정보

곡명, 아티스트, 가사, 멜론 DJ플레이리스트에 등록 되어있는 테마 데이터를 사용

데이터 수집

2) 사연 데이터 수집



가천대학교 대나무
숲
@gcubamboo

홈

정보

사진

동영상

게시물

커뮤니티

...



한양대학교 대나무
숲
@hyubamboo

홈

정보

사진

동영상

게시물

커뮤니티



4 졸업한지 좀 된 이별빠진 호랑이에요. 직장생활하면서 쉬는 시간이 생기면 고대
5 여러분 1월 1일은 잘 보내고 계신가요? 고대 학우분께 제 1월 1일을 자랑하고
6 (여성분들 답해주세요!) 안녕하세요. 남자입니다. 제가 어떤여성분을 알게되어서
7 대숲 한달쯤 전에 군대에서 이별했어요. 붙잡고 싶은데 붙잡으면 분명 돌아올까
8 2017년 1월 1일 0시 0분. 새해가 밝았습니다. 자정에 맞춰둔 알람이 울리고, 우리
9 신년을 맞아 변창하길 바라는 대숲!! 남사친 여자친때문에 논란이 많은데요! 남
10 당신은 항상 나이에 급급해한다. 뒤쳐지고 있는 것 같다고. 하지만 두려워 하지
11 지겨운 방학이 돌아왔다. 기말 전후로 소개팅을 통해 만났던 너 생각이 아무것도
12 여느 때와 없는 날들이다. 평소에 타던 33번 버스는 늘상 그랬듯이 우리 집 앞 정
13 새해에는..... 전역!!!!!! 전역을 하고 싶습니다!!!!!! 대한민국 군대 다 죽구
14 울 호랑이회님들,,,^^ 새해 복 많이 받으시구~^^하시시는 모든 일들도 변창하
15 이것은 유언장이 아니다 며칠에 걸친 지리한 싸움이 일단락됐다 어디서부터 잘
16 11시 45분. '이제 술술 올라가봐야하지않나?' 곧 도시 곳곳에서 피어날 불꽃들을
17 축하한다. 여러분은 벌써 올해의 1/365를 성공적으로 낭비했다. 대학생 (0/1/0)
18 즐거운 마음으로 참석한 새터(혹은 개강파티)... '나는 1X학번인데, 편하게 대해
19 몇년째 엄청 친한 여자친인데 한때 짝사랑의 대상이었지만 마음을 깔끔히 접었
20 어째서 교수님께서는 아직까지 성적을 공개하지 않으시는 걸까. 우리를 위한 나
21 고대 대나무숲 지기를 닉네를 ~냥으로 통일했던데..... 고양이인가 ㅇㅈ하는 부
22 노트북으로 드라마 보는데 화면 어두워질 때 비치는 내 모습이 존맛이라 볼 때
23 가끔 절 이유 없이 싫어하는 사람들이 있어요. 그러면 저는 그 사람들에게 이 악
24 정유년은 붉은 닭의 해입니다. 그렇습니다. 양념치킨의 해라는 거죠.
25 님들 솔직히 비뽀보다 그냥 에이가 더 기분 나쁘지않아요? 아 이번 학기 망했어
26 어둠은 빛을 이길 수 없다... 2017년 새해가 밝았습니다. 어둠은 조그마한 빛만
27 글을 잘 쓰고 싶다. 이과생인 내가 이런 재능을 탐하는 건 도통놈 심보라는 걸 알
28 나는, 간질임이다 여러분이 생각하는 그 간질은 아니에요 막 거품물고 쓰러져서
29 제가 군대 가는 것 때문에 여자친구가 너무너무너무너무 많은 걱정을 해서 고민

서울 시내 20개 대학에서 각 2000개 씩 크롤링



가천대학교 대나무
숲
@gcubamboo

홈
정보
사진
동영상
게시물
커뮤니티



고양대학교 대나무
숲
@nyubamboo

홈
정보
사진
동영상
게시물
커뮤니티

... 사연 데이터 약 40000만개 확보

4 졸업한지 좀 된 이빨빠진 호랑이예요. 직장생활하면서 쉬는 시간이 생기면 고대
5 여러분 1월 1일은 잘 보내고 계신가요? 고대 학우분께 제 1월 1일을 자랑하고
6 (여성분들 답해주세요!) 안녕하세요. 남자입니다. 제가 어떤여성분을 알게되어서
7 대숲 한달쯤 전에 군대에서 이별했어요. 붙잡고 싶은데 붙잡으면 분명 돌아올까
8 2017년 1월 1일 0시 0분. 새해가 밝았습니다. 자정에 맞춰둔 알람이 울리고, 우리
9 신년을 맞아 변창하길 바라는 대숲!! 남사친 여자친때문에 논란이 많은데요! 남
10 당신은 항상 나이에 급급해한다. 뒤쳐지고 있는 것 같다고. 하지만 두려워 하지
11 지겨운 방학이 돌아왔다. 기말 전후로 소개팅을 통해 만났던 너 생각에 아무것도
12 여느 때와 없는 날들이다. 평소에 타던 33번 버스는 늘상 그랬듯이 우리 집 앞 정
13 새해에는..... 전역!!!!!! 전역을 하고 싶습니다!!!!!!!!!! 대한민국 군대 다 죽구
14 을 호랑이원님덜,,, ^^ 새해 복 많이 받으시구~~^^ 하시시는 모든 일들도 변창하
15 이것은 우연장이 아니다 며칠에 걸친 지리한 싸움이 일단락됐다 어디서부터 잘
16 축하한다. 그분은 나 씨 올해 173cm 성공적으로 남비했다. 대학생 (0/1/0)
17 즐거운 마음으로 참석한 새터(혹은 개강파티)... '나는 1X학번인데, 편하게 대해
18 몇년째 엄청 친한 여자친인데 한때 작사랑의 대상이었지만 마음을 깔끔히 접었
19 어째서 교수님께서는 아직까지 성적을 공개하지 않으시는 걸까. 우리를 위한 나
20 고대 대나무숲 지기를 닉네를 ~냥으로 통일했던데..... 고양이인거 ㅇㅎ하는 부
21 노트북으로 드라마 보는데 화면 어두워질 때 비치는 내 모습이 존못이라 볼 때
22 가끔 절 이유 없이 싫어하는 사람들이 있어요. 그러면 저는 그 사람들에게 이 약
23 정유년은 붉은 닭의 해입니다. 그렇습니다. 양념치킨의 해라는 거죠.
24 님들 솔직히 비빔밥보다 그냥 에이가 더 기분 나쁘지않아요? 아 이번 학기 망했어
25 어둠은 빛을 이길 수 없다... 2017년 새해가 밝았습니다. 어둠은 조그마한 빛만
27 글을 잘 쓰고 싶다. 이과생인 내가 이런 재능을 탐하는 건 도동늬 심보라는 걸 알
28 나는, 간절입니다 여러분이 생각하는 그 간절은 아니예요 막 거품물고 쓰러져서
29 제가 군대 가는 것 때문에 여자친구가 너무너무너너너 많은 걱정을 해서 고민

서울 시내 20개 대학에서 각 2000개 씩 크롤링

진행과정

데이터 수집 ➤ 데이터 정제 ➤ 모델링 ➤ 추천시스템 구현 ➤ 결론

- 맞춤법 검사
- 불용어 처리, tokenize

으아아아아아아아아아아아아아아아 벋
 꽃 빨리 저버려!!!!!! 사람 마
 음 심송생송하게 하지말고!!!!
 꽃구경 가는 커플들 다 예쁜사람
 하세요!!!!!!!!!!!!!! 다만 내 눈
 앞에서만 사라져 줘요 제
 알!!!!!!!!!!!!!!!!!!!!

아;;;; 20 공부 안하고도 성적
 잘 받을 수 있다면 안암공전에서
 부러 저 멀리 교우회관까지 텀블
 링으로 투어하기 가능이에요 교수님
 ㅠㅈ... 공부시러요

으아아아아아아아아아아아아아아아 벅
 꽃 빨리 저버려!!!!!! 사람 마
 음 심송생송하게 하지말고!!!!
 꽃구경 가는 커플들 다 예쁜사람
 하세요!!!!!!!!!!!!!! 다만 내 눈
 앞에서만 사라져줘요 젠
 알!!!!!!!!!!!!!!!!!!!!

아;;;;; 20 공부 안하고도 성적
 잘 받을 수 있다면 안암공전에서
 부러 저 멀리 교우회관까지 텀블
 링으로 투어하기 가능이에요 교수님
 ㅠㅠ 공부시켜요

1) 맞춤법 검사



밀린 숙제를 제 시간에 끝마치지 못해서 허겁지겁 몰아서 했던 기억이 있다. 초등학교 때의 방학숙제부터 대학에서의 레포트까지, 밀린 숙제는 언제나 있었다. 밀린 숙제를 할 때마다 왜 이렇게까지 숙제를 미뤘었는지에 대한 후회, 과거의 나에 대한 원망이 엄하곤 했다. 그렇게 여러 궁시렁거림 속에서도 밀린 숙제는 언제나 끝이 있었다.

그리 오래 산 것도 아니지만, 너는 내게 밀린 숙제 같다. 그것도 조별과제. 조별과제이지.

413/500자 | 내용삭제

검사하기

밀린 숙제를 **제시간**에 끝마치지 못해서 허겁지겁 몰아서 했던 기억이 있다. 초등학교 때의 방학숙제부터 대학에서의 **리포트**까지, 밀린 숙제는 언제나 있었다. 밀린 숙제를 할 때마다 왜 이렇게까지 숙제를 미뤘는지에 대한 후회, 과거의 나에 대한 원망이 얹히곤 했다. 그렇게 여러 **구시형가름** 속에서도 밀린 숙제는 언제나 끝이 있었다.

검사결과가 복사되었습니다.

원하는 곳에 붙여넣기(Ctrl+V)해주세요

- 맞춤법
- 표준어익심
- 띄어쓰기
- 통계적교정



16



#2430번_제보 어제 관정 열람실에서 11시까지 공부하다가 깜빡하고 좌석 반납을 안했더니 바로 오늘부터 열흘 간 이용 정지 먹었어요.. 도서 반납도 30일까지는 연체가 가능한데 한 번 반납 안하고 집 갔다고 바로 이용 정지 맥이나요ㅠ 당장 중간고사가 코앞인데 너무 뻑이 치네요.. 도서관 이용하실 때 깜빡하지 맙시다..

유니코드, 문장부호, 숫자, 영어, 한자, 그 외 불용어 제거



"['어제', '관정', '열람', '실', '시', '공부', '깜빡', '좌석', '반납', '안', '오늘', '열흘', '간', '이용', '정지', '먹다', '도서', '반납', '까지는', '연체', '가능하다', '반납', '안', '집', '가다', '이용', '정지', '맥', '나오다', '중간고사', '코앞', '너무', '뻑', '치네다', '도서관', '이용', '깜빡', '맙시다']"

Tokenize 완료



#2430번_제보 어제 관정 열람실에서 11시까지 공부하다가 깜빡하고 좌석 반납을 안했더니 바로 오늘부터 열람 간 이용 정지 먹었어요.. 도서 반납도 30일까지는 연체가 가능한데 한 번 반납 안하고 집 갔다고 바로 이용 정지 맥이나요ㅠ 당장 중간고사가 코앞인데 너무 뻑이 치네요.. 도서관 이용하실 때 깜빡하지 맙시다..

유니코드, 문장부호, 숫자, 영어, 한자, 그 외 불용어 제거



총 1,600개

불용어 사전 오픈 데이터(1400개) + 직접 추가(200개)

"['어제', '관정', '열람', '실', '시', '공부', '깜빡', '좌석', '반납', '안', '오늘', '열람', '간', '이용', '정지', '먹다', '도서', '반납', '까지는', '연체', '가능하다', '반납', '안', '집', '가다', '이용', '정지', '맥', '나오다', '중간고사', '코앞', '너무', '뻑', '치네다', '도서관', '이용', '깜빡', '맙시다']"

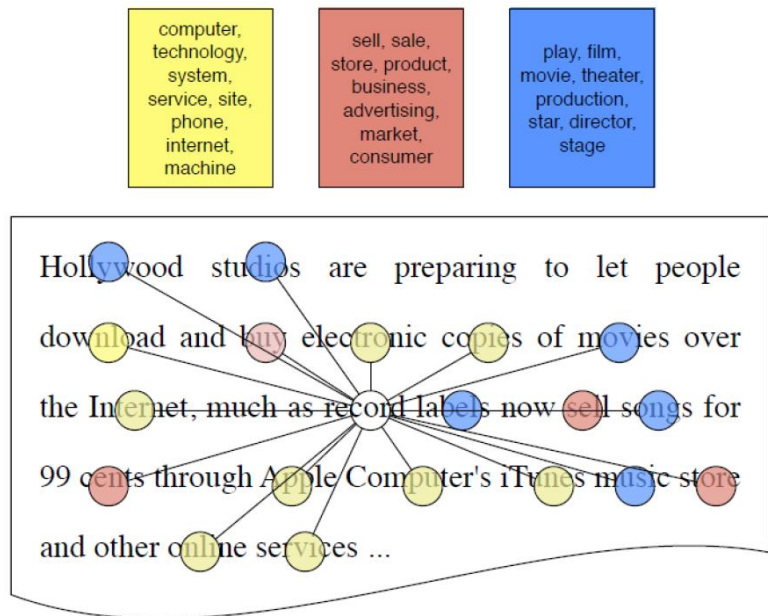
Tokenize 완료

진행과정

데이터 수집 ➤ 데이터 정제 ➤ 모델링 ➤ 추천시스템 구현 ➤ 결론

-토픽 모델링
-임베딩

1) LDA



사연 데이터

['어제', '관정', '열람', '실', '시', '공부', '깜빡', '좌석', '반납', '안', '오늘', '열흘', '간', '이용', '정지', '먹다', '도서', '반납', '까지는', '연체', '가능하다', '반납', '안', '집', '가다', '이용', '정지', '맥', '나오다', '중간고사', '코앞', '너무', '빡', '치네다', '도서관', '이용', '깜빡', '맙시']

음악 테마 데이터

100시리즈 2010년대 베스트 가요 감성적인
쓸쓸한 외로울때 발라드 발라드한 2018년
가을 송년회 늦가을 쌀쌀한 날 이별/슬픔
노래방 가요 가을밤 연도별음악 연도별
인기가요 2010년대 뮤직기네스

04 모델링

1) LDA - 사연

LDA에서 얻은 키워드 기반,
9개의 topic으로 분류

Manually
Naming

2+5: 일상/우울
4: 일상/밝은/학교
6: 새벽 감성
7+18: 일상/학교
8: 사랑/슬픔
9+11: 가족
10: 사회/답답함
16+19: 학교/우울
17: 사랑/설렘

Topic 수 = 20

개	적	억	돈	분	군대	아프다	학교	내	글	여성	아버지	오빠	노래	사람	애인	수업	안	가다	말다
두	인	밥	원	시	군	오늘	분들	사랑	엄마	수강신청	아버지	오빠	노래	사람	애인	수업	안	가다	말다
자리	되다	귀	소리	안녕하다	한국	눈	학생	사랑	아빠	정부	어머니	여동생	여행	기숙사	문	교수	너무	공부	술
이름	한	아침	게임	세상	전역	꿈	대	마음	없다	집단	이름	영화	세계	생각	열시	담배	친구	들다	선배
왔다	대해	방	인가요	입대	세상	밤	동아리	난	술다	우리나라	이름	영화	세계	보다	열다	냄새	친구	들다	제발
세	문제	를	팀	계시다	변만	밤	기간	만나다	대술	범피	이름	선물	화장	말	마르다	짜다	아니다	시험	당하다
성대	아니다	를	알머니	친해지다	군인	속	되다	다시	찾다	투표	주의	누르다	변화	못	내려가다	학부	그냥	오다	종강
휴학	대한	마시다	비	학번	하고도	술다	학우	날	대굴	검세	아름	카페	영상	되다	성매매	피우다	어떻다	열심히	인사
지키다	생각	저녁	놓다	후배	뉴스	손	많다	나르다	올리다	인권	사고	장	멋다	그렇다	침대	조	보다	다니다	씨
곳	않다	들리다	대체	새	실리다	길	생	이제	쓰다	회의	여름	주고받다	영여	심다	글자	강의	전	부모님	솔자리
중	위	맛있다	를	가요	복무	잠	드리다	헤어지다	싸다	복합	죽음	꽃	보상	더	색스	어색하다	진짜	대학	본인
지원	경우	음식	전	전공	팬	못다	질문	행복하다	있다	복합	폭력	크리스마스	시끄럽다	오르다	역	화장실	랑	학기	불편하다
번개	되이다	병	모으다	들다	인턴	하루	단	미안하다	올라오다	자유	무조건	잘해준다	불다	아니다	공연	토	들다	중	인간관계
실	성	기술이다	하	기업	눈물	정	나다	웃	키우다	정치	공정	꿈기다	전문	한	승위	왜냐하면	그렇다	정도	장난
확인	위해	맛	발굴	분위기	리	기억	활동	지내다	말다	대한민국	시선	별다	가사	지금	세다	술집	됐다	준비	유지
도서관	받다	참치	박동	입학	청춘	감사하다	동생	나라	당장	나라	드림	식다	출석	말이	가도	웃음소리	근대	생활	죽
이런	사회	바다	저다오다	터	두산	순간	관	서로	서강대	지치	죽이다	불법	탐	힘들다	유류	건물	애기	나간다	최종하다
사용	가지	심장	올리다	꼭	꼭소리	관	날기다	넌	페미니즘	연예인	프로필	만족하다	정말	큰일	구원	이제야	싫다	들어간다	싫다
알	때문	비용	꼭	잘만	꼭소리	관	날기다	넌	페미니즘	연예인	프로필	만족하다	정말	큰일	구원	이제야	싫다	들어간다	싫다
음악	어떻다	사주다	다리	배	나라	이별	가능하다	있다	국가	국가	통한	말다	그렇다	유류	걸	프로	현명하다	이렇다	조용하다
성취를	이유	들다	백	과목	외국	이경	떠오르다	아니	강정	통한	죄	드디어	태양	말다	무개	강의실	연락	늘다	조용하다
버스	의견	현드폰	과외	품	외국	외국	장학금	아직도	아직도	역지로	교내	목사	아픈거리다	이차	자다	스러운	대술	학년	끼리
덜다	이렇다	카드	통메	일로	작성	그녀	도움	지나다	감사하다	공지	가족	비밀	특징	해준다	범죄자	불	여자친구	면	누나
나무숲	개인	잔	저다오다	다녀오다	군무	지나가다	중	연애	매달리다	충고	바다	매달리다	매달리다	매달리다	매달리다	매달리다	매달리다	매달리다	매달리다
임	내용	끼	소문	성공관대	들다	예쁘다	겁니다	술프다	색	대표자	메아	싸	갑	알	살인	캠퍼스	남자친구	후	거들다
바꾸다	부분	힘내다	고양이	길	그다지	술	하나요	기억	겁	근거	장	이	성섭하다	이렇게	성관계	캐피드	오다	같다	식
공간	이라는	아어	신고	컴퓨터	배지	그리다	관련	처음	말리다	주의자	오른	말	관정	반다	전날	시점	건	고등학교	강강
선수	이해	시도	아저씨	신청	부대	행사	관심	연락	구체	장상스럼	일본	책	읽어버리다	말다	안동	덜	결다	살	뒤
부탁드리	일장	가게	연락처	내년	행사	행사	관심	연락	구체	장상스럼	일본	책	읽어버리다	말다	안동	덜	결다	살	뒤
들어다니	해결	영	양	화석	장점	별	관리	관리	마지막	종류	올림픽	식	위	사실	자제	감동	여자	반	달라
소	책	급하다	생리	여우다	직업	직다	필요하다	마지막	종류	올림픽	식	위	사실	자제	감동	여자	반	달라	달라

04 모델링

1) LDA - 음악 테마

LDA에서 얻은 키워드 기반,
9개의 topic으로 분류

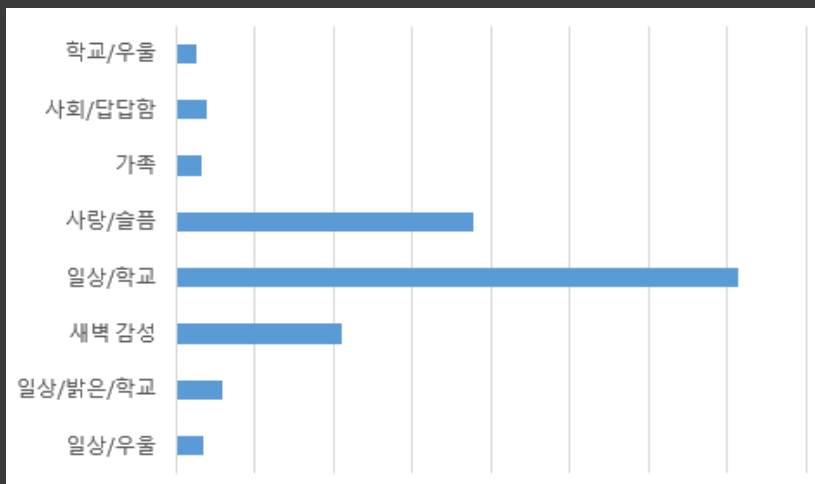
Topic 수 = 20

Manually
Naming

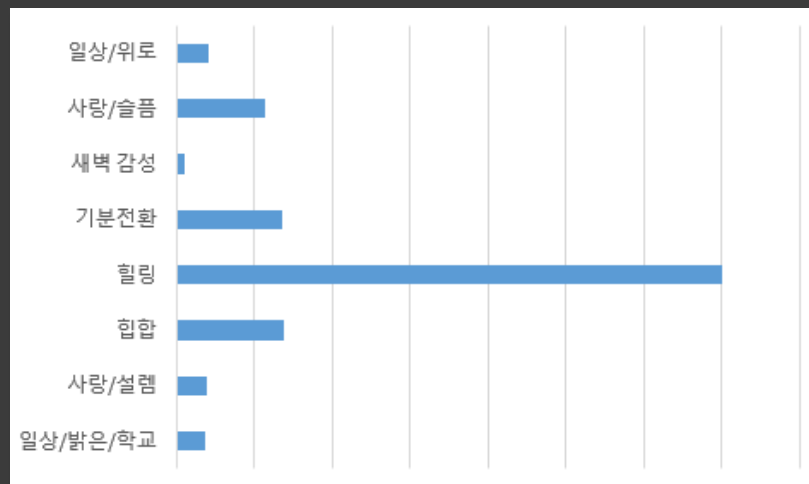
- 1: 일상/밝은
- 2: 사랑/설렘
- 3+17: 힙합
- 7: 힐링
- 8+12: 기분전환
- 11: 새벽 감성
- 14: 사랑/슬픔
- 15: 일상/위로
- 19: 가족

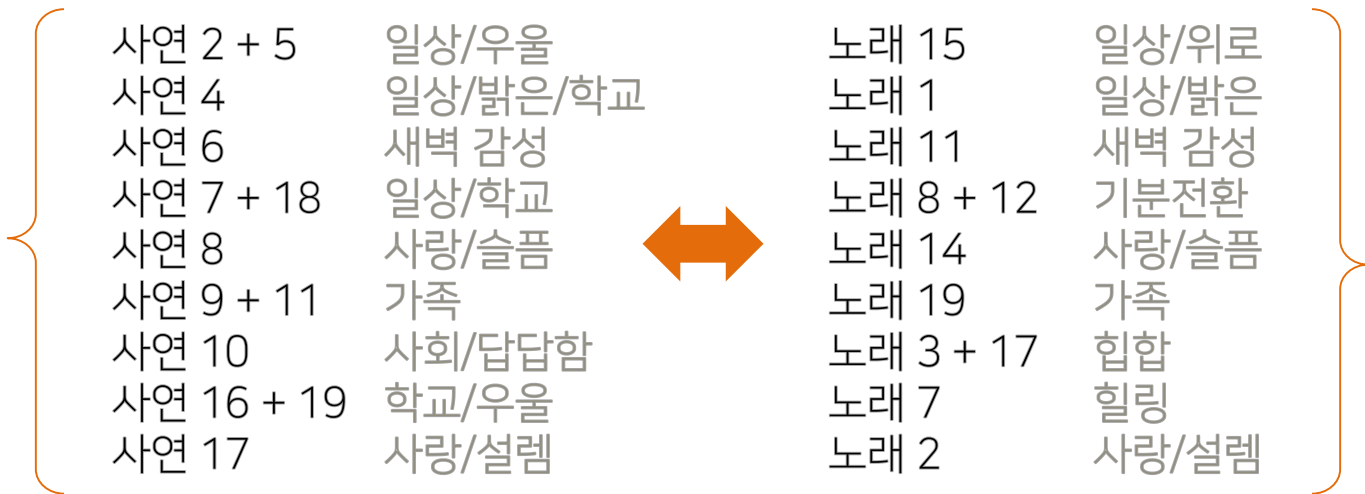
그녀	크다	함박함	연애	뉴잭스윙	꿈	봄	운영	안녕	유기	눈물	속	아침	비	우려하다	내	날	따르다	노래	은
여자	거기	직은	고독	별놓	세상	행복	영화	나이	야생마	있다	거짓말	빛꽃	내리다	관세	사랑	더	추다	라라라	억다
놀다	반	진진	별비	생마	하늘	일상	인연	커피	탕탕	아르다	진실	생	바다	백	천	추다	중	나라	집
얼마	달리다	빙그르	피자	다가다다	7위	키스	버스	노을	단물	기억	역	눈뜨다	여름	얼	었다	위로	음악	라라	꽃
가져가다	그림자	게리	빠라빠	밥	빛	주인공	저녁	연인	서울	떠나다	드라마	토	파도	갈매기	난	겨울	음	알라	눈
오오	맛	게릭	조죽	마야	소리	색깔	협합	정	날다	추여	안개	이마	아아	축하	말	인생	크리스마	렘	피다
오빠	죽	꽃점	메타	알리	만의	액주	따따	오후	방식	흐르다	추악하다	마리	비밀	부름	그대	고백	타	할라라	학
개	그걸	마이클	노노	빠빠	꾸다	불바람	헬로	장가	사	그림자	거짓	우와	시원하다	지구인	않다	그리울	즐거다	라할	지리
뻗다	불	호드도	미별	알다리	작다	데이트	세종대왕	분만	두두	이별	워워워	부	우	빙고	날	미련	히	멜로디	잔
아빠	올리다	이어온	나즈	뒹다리	뜨다	생애	파티	카페	산소	슬프다	달다	보아	신속하다	니	가을	이야기	빠르다	가을이다	출처
리	원	세노리따	코코	감성	뜨겁다	우비	장면	찰칵	네모	지우다	또다른	버려진	우산	정기석	망	가요	리듬	빠빠빠	배
노	머	강자별	노노노	뚜뚜뚜뚜	달다	팬	길가	케익	고향	을	가리우다	톡톡	루루	절다	잠	변만	키	우우우	맛있다
사료자다	세우다	더렸	알임다	똥주	새	해피	카카오	굿다	두르다	술	우우	빗	점배	심다	편지	아아	조물릿	자장	
반하다	군	이정표	켄나	롤러러러	구름	한강	돌대	집중	매매하다	떠나다	피란하늘	채리	바다로	점	오르다	전화가	소	슈가	다리
빠다	관심	아라	쿠라	끓는다	눈부시다	사랑	맛팔	지랄	쿠바	잊혀지다	서랍	리오	두번째	원장국	모르다	첫사랑	코	미가	술잔
찌릿	분	마카레나	밤밤	레베카	아래	마카롱	정류장	모금	아름	귀걸이	마하	그릴	두두두	스크림	말	그릴	다	술잔	금도
근대	감	내디디다	숙제	워고	달려가다	경쟁	불루스	로맨틱	이름	술름	물결	재다	엎다	평	사랑	낙엽	빨갳다	도비	검다
생일	룩	복일	코노	우연가진	날개	익	짜이다	엠	나머지	상처	외침	죽음	우연가	제자	오다	활	중추다	종알거리	발다
르	알으키다	특급	러즈	나오다	타고	돌쉴쉴	화이트	기차	자신감	남아	파도처럼	연날	바닷가	거리	한	온대	동네	내달	배달
까	에게도	호드	백죽	오라버니	발다	송	원색	공원	원색	흔들	흔들	나니	여름밤	상표	집착	알다	꽃노래	스르	
우후	산	드타	이키	스파이더맨	꿈꾸다	베이스	발표	샌드위치	아리랑	들어서다	오래되다	오노	그치다	세기	눈	가기	도록	재미	소개
물린	아지피	문위크	사라	쉬엄쉬엄	크게	겁니다	우연히	슈퍼	과객	미우다	웃이다	음식	섬	잘머니	결	연	올해	세레나데	텍스
를	켜	알싸다	아후	과파	날다	포옹	자유로	부노님	깃발	마음	흔든	미미	순수	마음껏	마음	말	잘생기다	마음껏	금도
남자친구	형	순정	우시	만주	놀다	호방	셀프	출리다	라인	지원하다	니나비	다미	클락	살라	초단	고단	파춥니다	무지	게임
아가씨	제로	트윈스	치카	나이팅게일	스르다	그랬는데	지영	비국	뒷모습	이니셜	지우	비처럼	허풍	넌	어젯밤	음지이다	슈비	스르	
완전하다	오케이	달로우	자옥	중화점	얼어나다	도란님	엠블럼	레퍼	매죽	와의	무대져	아우	푸른	산이	다시	소나기	비밀스런	년대	허하
오리	학교	락엔물	아시	가시리	오오오	장결	영화관	와인	주르르	시리다	술속	매롱	여인	외계인	가타	나인	나노	발차기	딱
참조	적당하다	구멍이	레노	스트레스	문	이모티콘	스납백	일요일	대한민국	어도	턱차	요리	루비	엘리베이	너무	하늘아	플릭	플루	부들다
소개	주말	우어	파파	롤러러러	눈빛	당배락	스나이퍼	복죽	공작	국배이	성당	즈	루비	장난감	와	점음	분위기	기막히다	달달

사연 Topic



음악 테마 Topic





Pair 별 embedding 공간 생성 (Doc2Vec)



Original Data

왜냐면 난 개를 좋아해 개가 어떤 앤지 알고서도 좋아해 개가 다른 남잘 만나는 걸 알고 있지 나는 너는 여신이니까
그래도 돼 너는 여신이니까 여신이니까 여신이니까 너는 되고 나는 다 안 되는 게 법이니까 난 그냥 널 봐서
강아지처럼 기분이 좋아 귀여워할 때만 나를 사용해도 좋아 호구 기리보이

Tokenized Data

"['왜냐면', '난', '개', '좋아하다', '개', '어떨다', '앤지', '알다', '좋아하다', '개', '남다', '만나다', '걸', '알', '여신', '이
니까', '돼다', '여신', '이니까', '여신', '이니까', '여신', '이니까', '되다', '안', '되다', '법', '이니까', '난', '그냥', '널',
'보다', '강아지', '기분', '좋다', '귀엽다', '나르다', '사용', '좋다', '호구', '기리보이']"



```
array([ 0.6298259, -0.09085082, -0.28519925, 0.5714372, 0.25872275,
        0.60149074, -0.02625632, 0.07765006, -0.7839311, -1.5680498,
       -0.92545795, 0.13430485, 0.36699116, 0.15147817, -0.96257496,
       -0.78277445, 0.5164371, 0.29842624, -0.205785, 0.22466424,
       -0.9545996, 0.6128009, -0.349194, -0.31204844, -0.7457065,
       -0.54703647, 0.1896811, 0.1491383, -0.21775162, -0.86559105,
       -0.48137563, -0.16984546, 0.8141323, 0.0564562, 1.0037832,
        0.81033397, -0.58311844, 0.5160924, 0.510567, -0.5937158,
        0.5981595, 0.87641585, -0.5774979, -0.40876126, 0.3795267,
       -0.4001868, 0.07840311, -0.7881859, 0.22949971, 0.42853683],
      dtype=float32)
```



Original Data

왜냐면 난 개를 좋아해 개가 어떤 앤지 알고서도 좋아해 개가 다른 남잘 만나는 걸 알고 있지 나는 너는 여신이니까
그래도 돼 너는 여신이니까 여신이니까 여신이니까 너는 되고 나는 다 안 되는 게 법이니까 난 그냥 널 봐서
강아지처럼 기분이 좋아 귀여워할 때만 나랑 사귀어도 좋아

총 80,000개의

Tokenized Data

사연, 가사 데이터에 대해

9개의 space에 embedding

"['왜냐면', '난', '개', '좋아한다', '개가', '어떤', '앤지', '알고서도', '좋아해', '개가', '다른', '남잘', '만나는', '걸', '알고', '있지', '나는', '너는', '여신이', '이', '니까', '돼다', '여신', '이니까', '여신', '이니까', '여신', '이니까', '되다', '난', '되다', '안', '되는', '게', '법이', '니까', '난', '그냥', '널', '보', '다', '강아지', '처럼', '기분', '좋아', '귀여워', '할', '때만', '나랑', '사귀', '어도', '좋아', '보이']"

```
array([ 0.6298259, -0.09085082, -0.28519925, 0.5714372, 0.25872275,
        0.60149074, -0.02625632, 0.07765006, -0.7839311, -1.5680498,
        -0.92545795, 0.13430485, 0.36699116, 0.15147817, -0.96257496,
        -0.78277445, 0.5164371, 0.29842624, -0.205785, 0.22466424,
        -0.9545996, 0.6128009, -0.349194, -0.31204844, -0.7457065,
        -0.54703647, 0.1896811, 0.1491383, -0.21775162, -0.86559105,
        -0.48137563, -0.16984546, 0.8141323, 0.0564562, 1.0037832,
        0.81033397, -0.58311844, 0.5160924, 0.510567, -0.5937158,
        0.5981595, 0.87641585, -0.5774979, -0.40876126, 0.3795267,
        -0.4001868, 0.07840311, -0.7881859, 0.22949971, 0.42853683],
      dtype=float32)
```

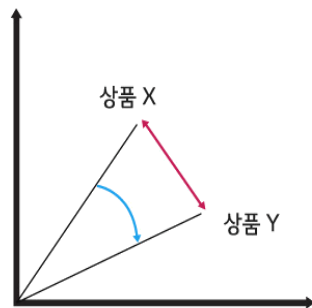
진행과정

데이터 수집 ➤ 데이터 정제 ➤ 모델링 ➤ 추천시스템 구현 ➤ 결론

-코사인 유사도 기반의
음악 5곡 추천

04 모델링

코사인 유사도 기반



← 유크리드 거리
↪ 코사인 유사도

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

사연과 유사도 높은 Top 5 음악

1

artist	title	lyrics
로시(Rothy)	다 핀 꽃	['핀', '이름', '모르다', '꽃', '송이', '떠나다', '전의', '글']...

산이(San E)	나 왜이래 (feat. 강민희 Of 미스에스)	['널', '바라보다', '표정', '관리', '안', '왜다', '얼굴', '써다']...
-----------	---------------------------	---

⋮

임정희	Golden Lady (feat. 현아 Of 4Minute)	['없이', '어떻다', '살아가다', '바보', '질문', '말', '알', '남']...
-----	-----------------------------------	---

먼데이 키즈 (Monday Kiz)	가을 안부	['어떻다', '자다', '지내다', '지난', '여름', '유난히', '힘', '...']
---------------------	-------	---

5

에디킴(Eddy Kim)	긴 밤이 오면	['기다', '밤', '오다', '길', '잃다', '니', '없다', '난하다', '...']
---------------	---------	---

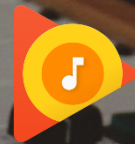
진행과정

데이터 수집 ➤ 데이터 정제 ➤ 모델링 ➤ 추천시스템 구현 ➤ 결론

05 결과



교수님이 어린이날에 보강을
한다네요. 휴일이라 친구들이랑
여행 계획도 짜고 예약까지
해봤는데 (...) 굳이 휴일에
수업을 해야하나요.



Play List

Title	Artist
Dangerous	보아 (BoA)
Pass	쥬얼리
A-G-E	E SENS
Tokyo Inn	혁오
Lonely(없구나)	B1A4

05 결과

교수님이 어린이날에 보강을 한다네요. 휴일이라 친구들이랑 여행 계획도 짜고
예약까지 해놨는데 (···) 굳이 휴일에 수업을 해야하나요.

추천 곡 1 (Dangerous- 보아)

... 장난치지 말고
속이려 들지 말고
서투른 변명 말고
Hurry up hurry up ...

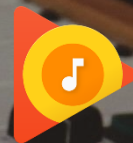
... 앞뒤 맞질 않지
넌 내손안에 있지
넌 참 뻔하게 보이는
습관적 거짓말 늘어놔...

추천 곡 2 (Pass- 주얼리)

... 이젠 PASS 너는 PASS
참을만큼 참았던 말야
답답해 지루해
늘 네 기준에 날 맞추고...



조현아 동영상을 보려고 했다.
 끝까지 보지 못하고 중간에 꺼야
 했다. 엄마가 소리 지르는 모습이
 오버랩 되었기 때문이다. (...)
 그렇다고 엄마가 날 사랑하지 않은
 건 아니다. (...)엄마가 없었으면
 지금의 나는 없었을 거다.



Play List

Title

Artist

엄마로 산다는 것은

보아 (BoA)

세노야

쥬얼리

별 헤는 밤

E SENS

그 아버지에 그 아들

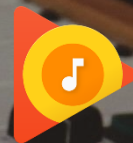
혁오

05

결과



선배 저한테 왜 그러셨어요.(...) 저
애정결핍 얼마나 심각한지
모르시잖아요. 문득 떠오르는 조각들이
저를 얼마나 힘들게 하는지
아시나요(...)그냥 마음 없던 척
선배 모르던 척 살아갈래요.



Play List

Title

Artist

헤어지기로 해

넬 (NELL)

Idiot

카라

연애 같은 걸 하니까

소란 (SORAN)

05 비교

1) 네이버 뮤직

교수님이 어린이날에 보강을 한
다네요. 휴일이라 친구들이랑 여
행 계획도 짜고 예약까지 해
놔는데 (...) 굳이 휴일에 수업
을 해야하나요.

검색 키워드

보강



Play List

TITLE

ARTIST

✓ 고급스런 명상 뉴에이지 일상의 휴식
(건강보강, 신경피로 완화)
✓ 늦은 후회 Viloet F

06 비교

2) 라디오 사연



대학 신입생입니다. 저는 시골
에 살아 소규모 학교를 다녔습
니다. 그래서 친구들끼리 우정도
깊고 친밀했습니다. (...) 친구
2명이 원하는 대학에 합격하지
못해 아직 공부 중입니다. 다시
모여서 재밌게 놀자! 파이팅!

로이킴, 정준영의 친한친구 사연 中



신청곡

슈퍼스타 - 이한철



Play List

TITLE

ARTIST



혜화동 (혹은 쌍문동)

박보람



zero

백예린

Summary

추천 LOGIC

- 1 대나무숲 사연 데이터 입력
- 2 데이터 정제: 맞춤법 검사, 불용어(stopword) 제거, Tokenization
- 3 LDA로 나눈 9개의 topic 중 가장 가까운 topic으로 분류
- 4 해당 topic과 matching 시켜 놓은 노래 카테고리의 임베딩 공간에서 Input data 임베딩
- 5 코사인 유사도가 높은 노래 5개를 output으로 추출

개선 사항

멜로디 정보를 Feature로 추가하면, 더 좋은 결과를 얻을 것으로 기대

Reference

Stopwords open data

- <https://github.com/stopwords-iso/stopwords-ko>
- <https://www.ranks.nl/stopwords/korean>
- http://blog.daum.net/_blog/BlogTypeView.do?blogid=0FyGV&articleno=7611081&categoryId=223515®dt=20080526063911



[감사합니다]