



유튜브 댓글 기반
마케팅 의사결정 지원 시스템
TrendPop

4팀: 이건모 | 고은지 | 신재솔 | 이지선

-
- 1. 팀원 소개
 - 2. 프로젝트 목적과 필요성
 - 3. 진행 과정
 - 4. 화면 구성
 - 5. 시스템 설계
 - 6. 선행 연구

목 차

- 7. 개발 과정
 - 7-1. 모델링 선정 | 라벨링
 - 7-2. 파인튜닝 과정 | 성능평가
 - 7-3. 키워드 추출
 - 7-4. 댓글 주제 분류
 - 7-5. Gemini 활용
 - 7-6. 이슈사항 및 해결 방안
 - 8. 시연
 - 9. 향후 계획과 기대효과
-

팀원 소개



눈꽃여왕 고은지
#모델링 #전처리
#파인튜닝



빠리시간 신재술
#크롤링 #키워드
#Gemini



갓생인생 이건모
#프론트 #전처리
#파인튜닝



고양이손 이지선
#기획 #UI 설계
#QA

프로젝트 목표

[연예공화국]<5>전 세계 한류 침범 BTS→블랙핑크→뉴진스

관성훈 기자 odrom@imaeil.com

매일신문 입력 2023-12-16 07:00:00 수정 2023-12-17 18:12:29

5년 만에 글로벌 문화 영향력 7위에 입성, 24단계 도약

‘불어라 한류’ 연간 생산유발효과 10조원 안팎

LPGA 황금세대 ‘박세리 키즈’ 이후처럼 한류도 위기 올 수 있어

대한민국은 연예 강국이다. 전 국민이 연예인(셀럽)에 열광하고, 어릴 때부터 꿈이 대다수 '연예인'이다.

세계 문화대국 한국의 위상이 갈수록 높아지고 있다. 글로벌 문화 영향력은 2017년 세계 31위에서 2021년 7위로 도약했다. 이제는 미국과 유럽 문화강대국(영국, 프랑스, 이탈리아, 스페인) 그리고 이웃나라 일본과도 어깨를 나란히 할 정도로 소프트파워를 갖게 된 것이다.



프로젝트 목표

[연예공화국]<5>전 세계 한류 침범 BTS→블랙핑크→뉴진스

관성훈 기자 cdrom@imaeil.com

매일신문 입력 2023-12-16 07:00:00 수정 2023-12-17 18:12:29

5년 만에 글로벌 문화 영향력 7위에 입성, 24단계 도약

‘불어라 한류’ 연간 생산유발효과 10조원 안팎

LPGA 황금세대 ‘박세리 키즈’ 이후처럼 한류도 위기 올 수 있어

대한민국은 연예 강국이다. 전 국민이 연예인(셀럽)에 열광하고, 어릴 때부터 꿈이 대다수 ‘연예인’이다.

세계 문화대국 한국의 위상이 갈수록 높아지고 있다. 글로벌 문화 영향력은 2017년 세계 31위에서 2021년 7위로 도약했다. 이제는 미국과 유럽 문화강대국(영국, 프랑스, 이탈리아, 스페인) 그리고 이웃나라 일본과도 어깨를 나란히 할 정도로 소프트파워를 갖게 된 것이다.

니네 바이블로 뜬줄 알고 있었는데 이번에 사실을 알게 됐다.. 무슨일이 있어도 응원한다. 이번 사건을 끝까지 기억하고, 그 어떤 사이비에도 물들지말고 굳세게 살아라

👍 658 💬 답글

✓ 답글 7개

NewJeans (뉴진스) 'Bubble Gum' Official MV

HYBE LABELS

조회수 1,397,340

댓글 108,131개

댓글 추가...

@홍영웅-3kg 5개월 전
0:55 도서의 제목은 <순수>의 시대>, 이 책에서 주인공은 관습에 따르면서 수감생활이 암시적일거 같습니다

@이르케를게 5개월 전(수정됨)
책인가 비노발로 만드는데 아르케를게 너무 고마워!— 또 올게

@99226 5개월 전
아무리 봐봐도 미아도 뉴진스가 진짜다... 나조차도 다른 건 그냥 잊혀진 그냥 뉴진스는 뉴진스야

@py1047 1개월 전
간만에 토이 Don't be blue!

@whosthatman 5개월 전
불어 아연양 뉴진스랑 끝내주는 3월보내고 개학하기 전날밤잠들

@홍우미용단-kbr 5개월 전

@the_age_of_romance 11일 전
나는 너희들의 유자구름도 깨지지 않았으면 좋겠어

그냥 뉴진스는 시대의 상징임

독보적인 존재라고 생각함

👍 5.3천 💬 답글

✓ 답글 21개

프로젝트 목표

[연예공화국]<5>전 세계 한류 침범 BTS→블랙핑크

권성훈 기자 cdrom@maeil.com

매일신문 일백 2023-12-16 07:00:00 수정 2023-12-17 18:12:29

5년 만에 글로벌 문화 영향력 7위에 입성, 24단계 도약

‘불어라 한류’ 연간 생산유발효과 10조원 안팎

LPGA 황금세대 ‘박세리 키즈’ 이후처럼 한류도 위기 올 수 있어

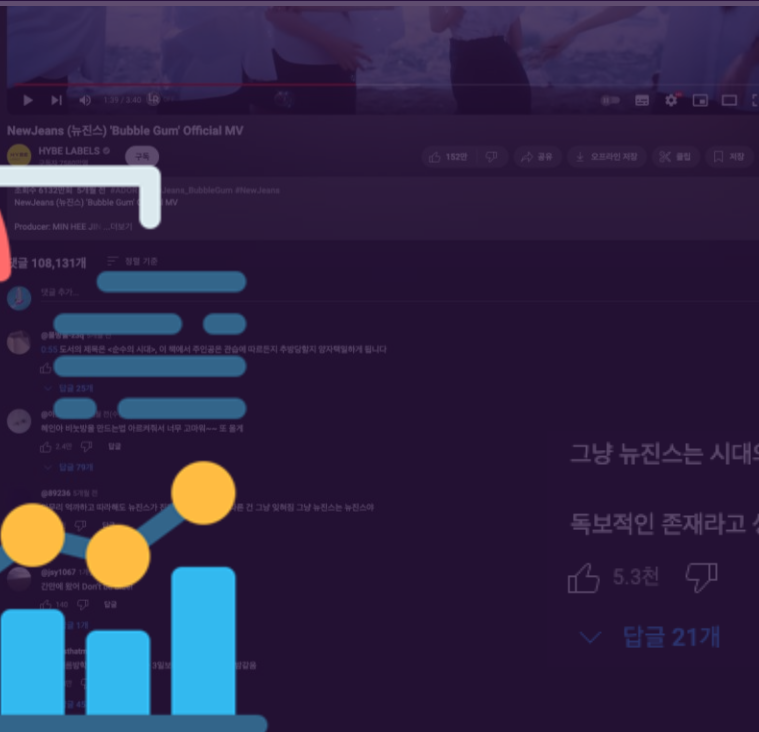
대한민국은 연예 강국이다. 전 국민이 연예인(셀럽)에 열광하고, 아랍 페트리 꿈이 대다수 ‘연예인’이다.

세계 문화대국 한국의 위상이 갈수록 높아지고 있다. 글로벌 문화 영향력은 2017년 세계 31위에서 2021년 7위로 도약했다. 이제는 미국과 유럽 문화강대국(영국, 프랑스, 이탈리아, 스페인) 그리고 이웃나라 일본과도 어깨를 나란히 할 정도로 소프트웨어를 갖게 된 것이다.

니네 바이블로 뜬줄 알고 있었는데 이번에 사실을 알게 됐다.. 무슨일이 있어도 응원한다. 이번 사건을 끝까지 기억하고, 그 어떤 사이비에도 물들지말고 굳세게 살아라

658 답글

답글 7개



그냥 뉴진스는 시대의 상징임

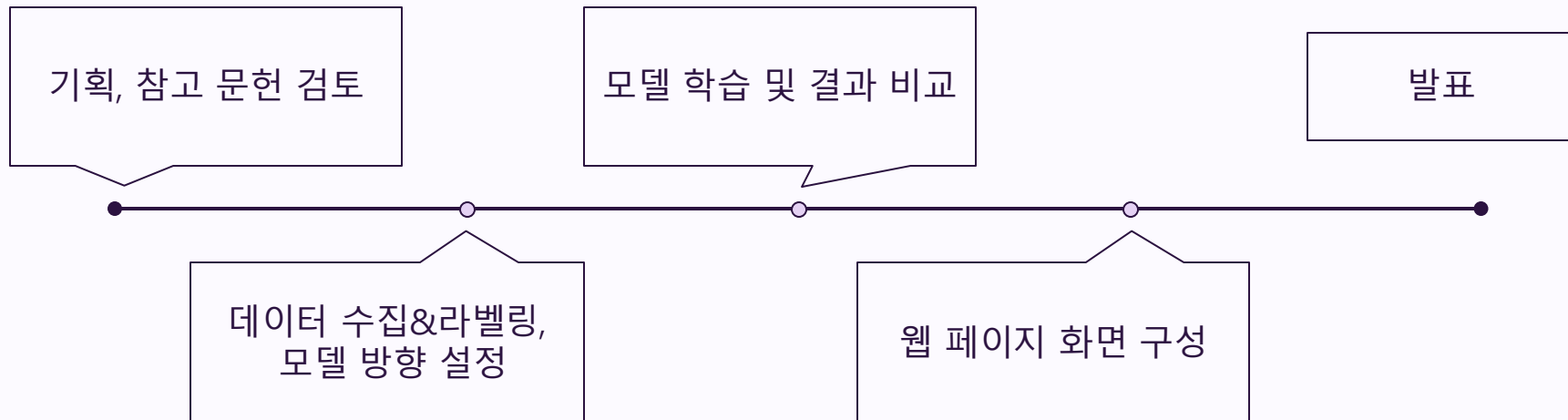
독보적인 존재라고 생각함

5.3천 답글

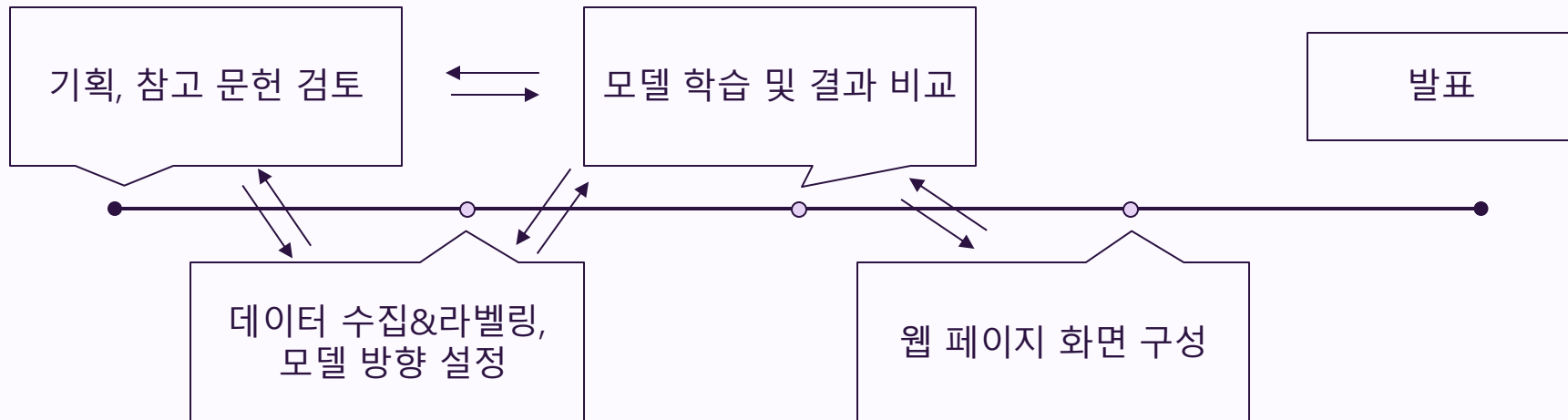
답글 21개

@the_age_of_romance 11월 31
나는 너희들의 유자구출을 위해서 있었지만 불행이
25 208 12 답글

진행 과정(예상)



진행 과정(실제)



화면 구성



Top 10 아티스트 그룹 선택

1 전체 댓글 반응 감성 분석
차트

2 Gemini로 뉴스 요약 검색(RAG)

뮤직비디오 선택

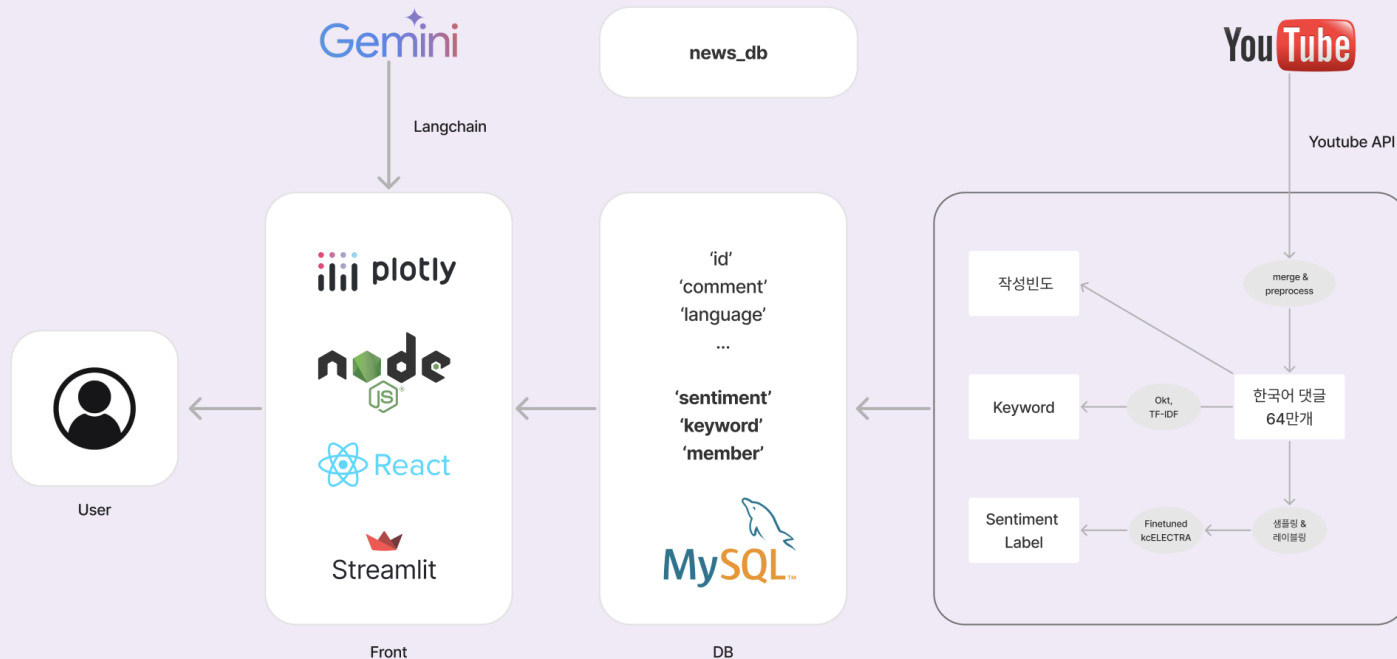
3 뮤직비디오 별 반응 감성 분석
차트

멤버 선택

4 그룹/멤버 별 키워드 분석
기능

5 키워드, 감성으로 댓글
파라미터

시스템 설계



선행 연구

제목	연도	데이터/레이블링	접근	전처리	모델	학습/성능향상	성능평가	
			어휘기반(Lexicon), ML기반(ML), 신경망기반(NN), LLM API(Prompt)	형태소분석기, 불용어처리, 명	감정점수합산, SVM/Boostin	비지도학습/사전학습, 지도학		
2013 KOSAC - 한국어 감정 및 의견 분석 코퍼스	2013		Lexicon	형태소분석	svm, 군부정 극성, 주관적(se		10류를 교차 검증법 사용	
2019 자연어 처리 기술을 이용한 감정분석 기법에 관한 연구	2019	트위터데이터 크롤링	Lexicon	형태소분석기	비학습데이터, 기계학습 병행			
2020 기계학습을 이용한 Aspect-Based Sentiment Analysis 기반 전가차	2020		ML, NN	불필요한 텍스트 요소 제거 (U 토큰화 및 품사 태깅 표제어 추출(Lemmatization) 불용어(Stopwords) 제거 데이터 불균형 해결을 위한 오	요소 추출: TextRank, Naïve 감성분석: Logistic Regressi 임베딩 : TF-IDF, GloVe			
2024 한국어 구문분석을 활용한 의존 관계 패턴 기반의 감성사전 구축 기법	2024	소평문 후기 분석 데이터	Lexicon	KSS 문장분리기, 한 어절 이	의존 구문분석, 규칙 기반 접근	비지도 or 지도 X 의존 관계 학습	F1 score	
2019 BERT 기반 한국어 감정 사전을 이용한 감정 예측기 개발	2019	온라인 댓글 및 리뷰데이터 수	NN	토큰화, 형태소분석기 - 불용어	nn사용	nn사전학습, 파인튜닝, 하이퍼: 정확도, f1스코어, 정밀도, 재현율		
2023 초거대 언어 모델을 활용한 감정 분석 연구	2023	GoEmotions(레딧 댓글) 'K	NN, Prompt	형태소분석기, 동사 형용사 추	BERT(KcBERT, RoBERTa GPT(3.5, 4)	사전학습, Zero-shot prompt	F1 score, Pearson 상관계수	
2016 한국어 감정분석 코퍼스를 활용한 양성정보 기반의 감정분석 연구	신호필, 김	2016		표현유형 태그, seed태그추적	선형 svm,scikit-learn	확률자질 (PES, PPOSS, P REF-TYPE-SPEECH, REI	10겹 교차검증법, precision, recall, f1, accuracy	
LLM과 양상별 머신러닝모델을 융합한 한국인 감정분석 모델 설계	김현지	2024	AIHub 감정 대화 데이터셋 레이블 - 기쁨, 슬픔, 분노, 불	ML + NN	Konlpy.Okt 명사추출 FastText를 통해 각 '단어'를	KoBERT + RF 양상별	FastText - 트레이닝 KoBERT - 트레이닝	Accuracy, Precision, Recall, F1-Score
2021 Building the Korean Sentiment Lexicon for Finance (KOSELF) 조수지	2021	한국어금융감성사전, 애널리스트	lexicon	토큰화, konlpy형태소분석기,			로지스틱회귀분석, 회귀분석	
리뷰 감성분석 - KoBERT, KoGPT-2, KoBART	이민아	2023	구글 플레이스토어 카카오톡/ 평점 1~2 부정, 3 요청, 4~5 ; NN	5자 이 하 제거, 한글 자/모음단 고정된 입력 크기를 유지(64 t	KoBERT, KoGPT-2, KoBA	KoBERT, KoGPT-2, KoBA 하이퍼파라미터 튜닝 - 배치사	Accuracy, Precision, Recall, F1-Score	
라디오 청취자 문자 사연을 활용한 KoBERT 기반 한국어 다중 감정 분석 연구 이재아	2023	1) AIHub 감정 대화 데이터셋 2) 라디오 청취자 문자 사연을 7가지 - 행복', '슬픔', '놀람', '기	NN	15,831문장 -> 5,531문장 비문법적인 노이즈를 가진 문? '내용 없음(사진)', 중복된 사연	KoBERT	KoBERT 파인튜닝 모델1 - AIHub 데이터로 학습 모델2 - 라디오 청취자 데이터	Accuracy(라디오 청취자 test set)	
[감정 분석 모델] 한국어 감정 분석 데이터셋 KOTE 논문 리뷰 / Python에서 K	2024	KOTE 데이터셋	ML	최소 길이는 10, 최대 길이는 : 가장 긴 텍스트의 상위 10%의 길이는 404, 평균은 57.32, 종	KoELECTRA 를 KOTE로 파			
뉴스 기사 제목의 감정 라벨링 기법 연구	하재룡	2022	코로나19 관련 뉴스 기사 제목 Naver 포털에서 월별로 약 20 공정, 부정 수작업 레이블링	관용구, 특수문자 등 제거	1) KNU Lexicon 2) KoELECTRA	모델1 - 도메인 단어 사전을 거 모델2 - KoELECTRA 파인튜 모델3 - 추가 파인튜닝(NSMC 모델4 - 세가지 모델을 적절한	Accuracy, Precision, Recall, F1-Score	
딥러닝을 활용한 감정 분석 과정에서 필요한 데이터 전처리 및 형태 변형	서해진, 손	2022	0) EDA 1) 소문자화(영어) 2) 토큰화 3) 축약어 4) 태그 제거 5) 이모티콘 제거 6) 외국문자 제거 7) 무의미하게 반복된 문자 제거 8) 명사추출 9) 불용어 제거 9) 너무 긴 문장 제거					
2019 자연어 처리 기술을 이용한 감정분석 기법에 관한 연구	윤태성	2019	현대 한국어의 어휘빈도 자료	형태소분석기, konlpy.org,				

모델링 선정 | 라벨링



LLM labeling



KcELECTRA
Finetuned by
Kotedataset



Labeling



Fine-tuning

Group: NewJeans

Title: OMG

울분이 터져나온다 나라를 두번 잃을수는 없다

Sentiment:

List

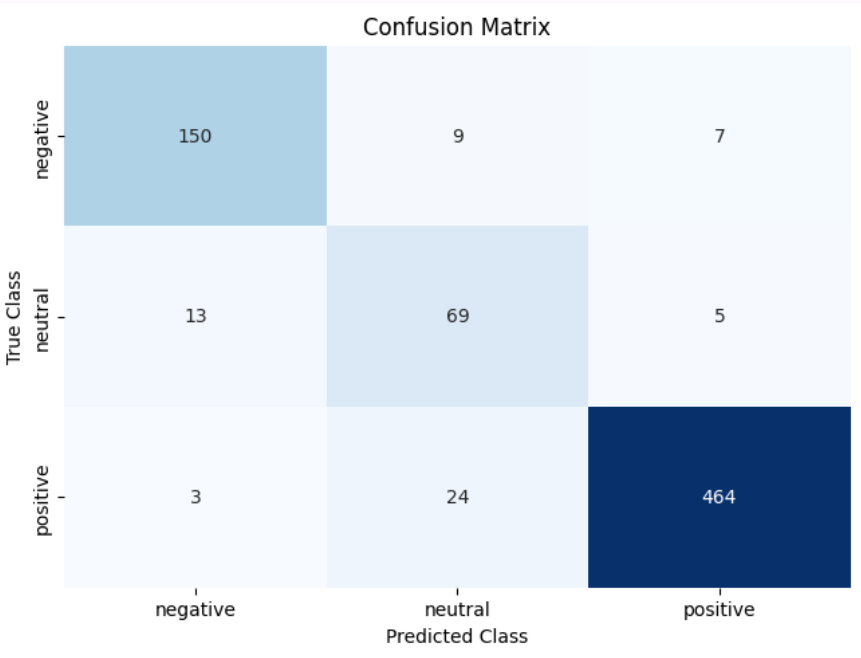
종료

직원들아 탈 하이브는 지능순..
타돌팬이지만 객적으로 잘됐으면하는 그룹 다 이쁘고 노래도 잘하고 노래가 잘 어울리네요 듣기편하고 화이팅!
놀아줘클럽헤인입니다
난지금까지는DIVE아니였지만 지금부터DIVE될거같아
신곡 미쳤네ㅠㅠㅠ
와미친...
역시 독보적 뉴진스
오늘 주주총회는 열리지 않습니다. 오늘 민회진이 위험하다고 퍼뜨리면 안됩니다. 역공당할 수 있어요. 물론 사이버 단체의 비
뉴진스 안좋아했는데 방시혁때문에 뉴진스 팬됨ㅋㅋㅋ
노래 진짜 구리다
→ 울분이 터져나온다 나라를 두번 잃을수는 없다
노래 진짜 넘좋자나 ㅠㅠ
르세라핌 노래들은 전부 다 명곡이야.. 더 대박나자아
흠흠. 닝닝을 매우 좋아한다
[timecode] 언니.....

파인 튜닝 | 성능 평가

Performance Metrics by various models

Model	Accuracy	Precision	Recall	F1-Score
KcELECTRA 12000 sample finetuned	0.918011	0.851625	0.880576	0.864481
gpt-4o with prompt	0.905914	0.824174	0.841612	0.832502
gpt-4o	0.90457	0.829732	0.866307	0.844036
KcELECTRA 17000- oversampled finetuned	0.895161	0.820019	0.852437	0.833488
gpt-4 finetuned with 283 examples	0.889785	0.824535	0.741701	0.768945
koELECTRA(Max)	0.880376	0.800417	0.781084	0.790099
gpt-4o-mini with prompt	0.879032	0.781747	0.810755	0.793331
koELECTRA base model	0.876344	0.840351	0.723631	0.749794
gpt-4o-mini	0.86828	0.780447	0.794009	0.777354
kote-ELECTRA 200- finetuned	0.643817	0.399864	0.413798	0.37505



Member Keyword Table

키워드 추출

닉네임 사전기반 멤버명 인식



키워드에서 타
멤버 언급제거

키워드 추출 : Count, TF-IDF

댓글 데이터

명사추출 : Okt와 Mecab의 비교

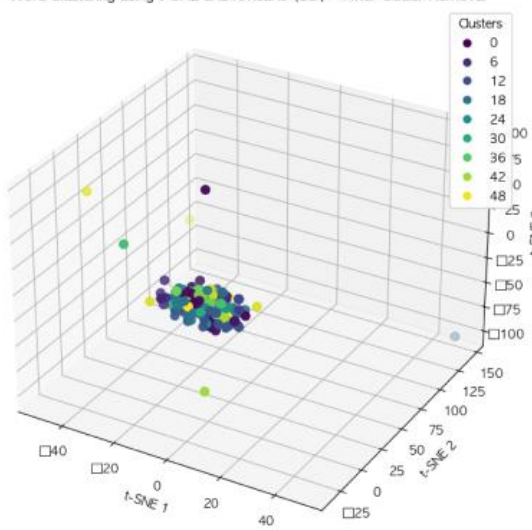
키워드 추출 : Count, TF-IDF

넷글 주제분류

=== 각 토릭에서 중요한 단어 5개 ===

토릭 1: 0.201*“사랑” + 0.052*“별로” + 0.044*“응원” + 0.028*“항상” + 0.019*“의상”
 토릭 2: 0.484*“믹스” + 0.034*“와이팅” + 0.019*“활동” + 0.015*“이제” + 0.012*“지
 토릭 3: 0.074*“릴리” + 0.050*“파트” + 0.046*“설문” + 0.042*“규진” + 0.039*“지
 토릭 4: 0.052*“대박” + 0.034*“그룹” + 0.027*“하나” + 0.017*“중간” + 0.016*“한
 토릭 5: 0.297*“노래” + 0.039*“유비” + 0.030*“이벤” + 0.027*“느낌” + 0.027*“중
 토릭 6: 0.033*“그것” + 0.029*“보컬” + 0.028*“타이틀곡” + 0.027*“수목” + 0.025*

Word Clustering using t-SNE and KMeans (3D) - After Outlier Removal



Clustering

장 관련 있는 토릭 ===
 (노래), 확률: 0.4087
 (스타일), 확률: 0.3717
 (사랑), 확률: 0.2831
 (노래), 확률: 0.5500

엔믹스 노래 중독성 있고 좋은데 왜 안뜨는지 진심 의문
 대쉬랑 비슷한것 같기도하고 이런 좀 다른 대성성있는 노래로갔으면하는 아쉬움이네요 ...

분부장이 피씨로 열병 클릭하면 별 부비로 연결 안되고 수정별인가 그 장면 있는 곳에...

이야 노래 좋아 맛있네

외모핵을 필두로 더 행방조길러다로서 해원이 얼마나 노력했을까 멤버들 모두 착하고 이...

노래는 좋은데 가사는 개짜지는 별

해원이 외모핵 이후로 전편 멤버 모두 착하고 예뻐 실력도 탑이네

와 설문이용 지우랑 비주얼 폭발하는 노래 하나 타이틀로 내자 엔믹 요즘 솔을...

이러온 까고 들으니가 무슨 연극같음

구름

이 노래가 위 라고위에서와 안에 들어갈을 알았는데

노래가 연이따나 대형기획사에서 낸 유비 조희수가 준수 길그룹보다 못하네 강 애들...

사랑해요

규진이 고정웃음 왜이리이지

중간에 엔믹스 나오는 부분중 그간행으면 좋겠어 맥골거

소녀감성 다 필요없어영생 합함해위

밥 아티스트 최초 빌보드 라틴 뮤직 위크 차트 유류보로 직접 검색 언제나 반전매력과...

규진 설은 이브당

솔직히 실력으로는 엔믹스가 원탑 아니나

대박 대박 안 말하고 실은데도 막 나오네요

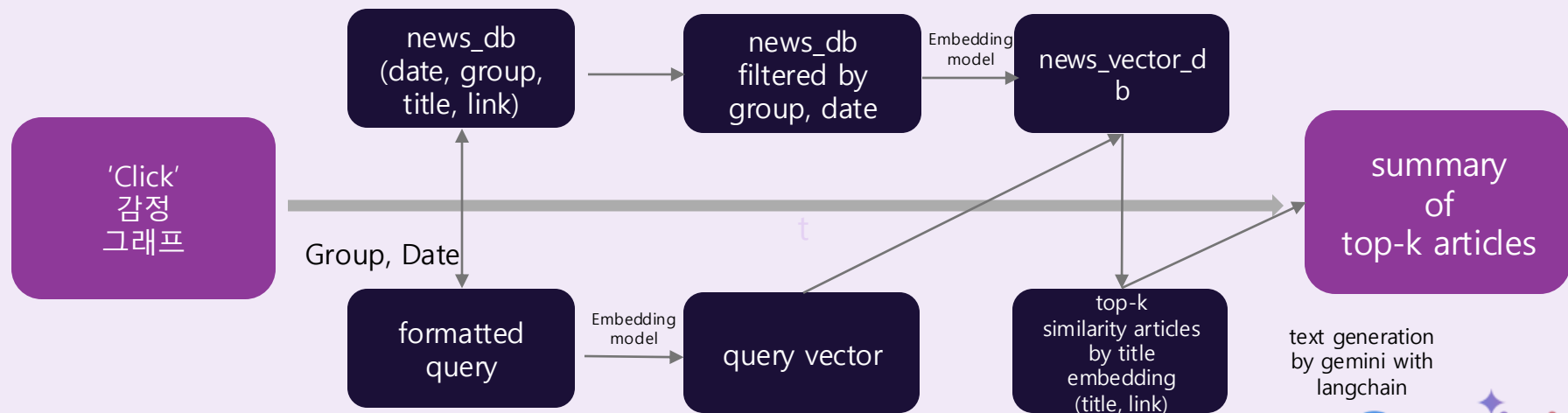
베이스 미침

G	I
comment	키워드
와씨 이번 한박곡 미쳤다 (아직 안봄) 보고 왔는데 개지리누 출연송 영상 빨리 빨리 올라와라!!	컴백, 댄스, 응원
오늘만 기다렸다	컴백
이 노래들이 정말 그냥 대중적으로 가자 제발 음악성 어떡고자짜고 모르고듣고 듣기 편하고 신나는걸로 가자	음악, 컨셉
근데 뉴진스 원래 조희수는 잘 안나올 대위곡부터도 조희수 일리는데	조희수, 순위
파종나네. 하필 저기로 노래 부르다니..	응원
와 미쳤다(아직안봄)	응원
요즘 지우 설은 많이 일어나주는 것 같아서 좋은 것같다고 기존에 일어했던 멤버들도 다 맞는 퍼트 주고.. 요즘 멤버, 퍼트, 소속사	1위 가조아
1위 가조아	응원, 순위
새창마지들을 상대하거면 많이 부족하다 너무 의식해서 오마하지 않길 바람	순위, 컨셉, 타그
엔믹스스리위서 좋네 o.o랑 다이아 개좋아해서...규진님은 진짜 미쳤다	음악, 컨셉
하루스름 케이팝은 언제나 올라....	음악, 장르, 하우스
니 인생이 훨씬 망함	키보드배틀
화이팅	응원
the song is really good keep up the enthusiasm for nmixx. i know nmixx will be more success! 기타	기타
KYUJIN 사 랑 해 요 !!!!!!!!!!!!!!!	응원
안보고싶으면 보질마세요 *	키보드배틀
빌보드 핫 트렌딩 송 1위 추카추카	응원, 순위, 빌보드
NMIXX Like OOH-AHH(OOH-AHH하게)	음악
아니 진심으로 엔비로서 회사서 그림 먹고 미친거임?? 4세대 여돌만 거의 전편터엔데 진심으로 이런걸로 송	음악, 속사, 컨셉
이게 엔믹스다	응원, 컨셉
사실 내심 이번엔 대중적으로 오나 기대했는데 한변 들으니가 대중적이고 뭐고 우린 이게 맞는 거같다ㅋㅋㅋ	음악, 컨셉
원 노래만가 했더니 제품달 ost 있네	음악
노래 맛있게 진짜...이게 엔믹스지	음악, 컨셉
지우 진짜 이쁨!!!!	멤버
와어 3분 왜이렇게 짧아요	음악
인근을 1위까지 약 백만뷰 남음~~ vamos!!!!!!!!!!!!	조희수
NMIXX (J)'s Top Popular Songs on YouTube Music! 별별별(See That?) - 35M DASH - 84M Love 기타	기타
이번노래도 대박 나자	음악, 응원
아.. 대쉬로 고정하고 이번엔 신장 둘을 총 알았더니 여기서 브레이크 걸리네... 많이 아쉬운데... 또다른 곡	음악, 컨셉
우리 딸을 너무 멋지고 재능이 많아서 매번 놀랍니다!!	응원
이게 노래라고?..	음악, 컨셉, 비난

사전할당 주제어(대주제)
 분류
 *향후과제

Topic Modeling(LDA)

Gemini 활용 이슈요약봇



text generation
by gemini with
langchain



와이어프레임

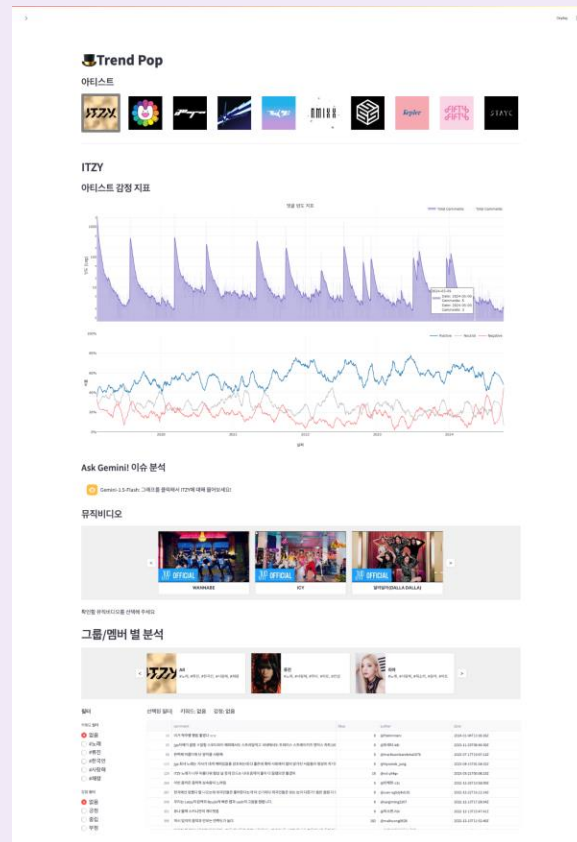


시연

와이어프레임



구현 완료



이슈 사항 및 해결 방안

실시간 대응
능력 향상

캠페인
성과의
즉각적 평가

트렌드
선도와
맞춤형
콘텐츠 제공

악플러 대응
및 리스크
관리

팬덤과의
소통 강화

향후 계획과 기대효과

[감정분석]

감정분석 모델 고도화

한국어 모델 성능개선

하이퍼파라미터 튜닝
lr, batch_size, n_epochs,
decay, dropout,
optimizer...

보다 많은 레이블링

(실험적) RAG를 통한
맥락 반영

유튜브 댓글의 특성과
입장, 감정, 맥락에 대한
보다 깊이 있는 연구

[다국어지원]

영어기반모델
(+Translate API)
다국어지원모델

[Realtime]

DB 업데이트 자동화
실시간/배치방식으로
댓글 및 뉴스 수집

[키워드]

키워드 방식보다 더
나은
주제분류 모델 고안

주제분류 지도학습을
위한
레이블 생성

댓글의 특성을 더 잘
반영하는
임베딩 모델 학습

[요약봇]

검색(임베딩) 정확도
개선
속도개선
스트림방식 적용

[웹호스팅 및 서비스]

GCP 및 Vertex AI를
통한
GCP Native Stack
구현

Compute Engine,
Storage,
Vertex AI,
BigQuery

감사합니다.
THANK YOU.
WISH YOU ALL THE BEST!

심사위원님들 피드백

-학습한 데이터 출처

-감정분석 영화 리뷰 데이터를 기준으로 한 이유? 비슷하지 않다

-데이터 전처리 과정을 어떻게 했는가?

-차트에 포커스 표시와 X축 기간 레이블

-감성 사전 만드는 노력!그게 장표에 없다!했던 노력에 대한 흔적을
녹여내라