

# 현대대중가요 가사 분석

172114 임성애

## 목차

1. 프로젝트 개요

2. 프로젝트 진행 과정

3. 프로젝트 결론

## 1. 프로젝트 개요

현대대중가요 장르별&실시간 분석을 통해 대중의 생각과 관심을 빠르게 파악

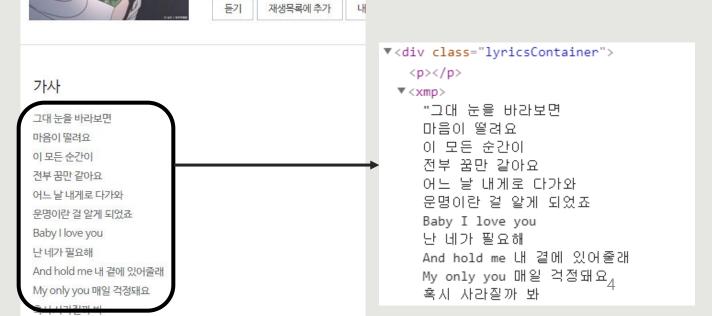
빅데이터 시스템을 이용한 국내 음원 '장르별 Top 100 가사' 분석

Top 100 가사 가져오기



ex. 발라드 Top 100





Top 100 가사 가져오기



['그대 눈을 바라보면\r\n마음이 떨려요\r\n이 모든 순간이\r\n전부 꿈만 같아요\r\n어 느 날 내게로 다가와\r\n운명이란 걸 알게 되었죠\r\nBaby I love you\r\n난 네가 필 요해\r\nAnd hold me 내 곁에 있어줄래\r\nMy only you 매일 걱정돼요\r\n혹시 사라 질까 봐\r\n그댄 나의 사랑\r\nOnly you\r\n내게 보였던 그 미소\r\n다 너무 예뻐서 \r\n더 보고 싶어서\r\n나를 달래고 있죠\r\n한 번에 알아 본 사랑은\r\n내 생에 오직 그대뿐이죠\r\nBaby I love you\r\n난 네가 필요해\r\nAnd hold me 내 곁에 있어 줄래\r\nMy only you 매일 걱정돼요\r\n혹시 사라질까 봐\r\n그댄 나의 사랑 \r\nOnly you\r\n푸르고 짙은 계절이 찾아올 때\r\n그대 나의 꽃이 되어 주기를 \r\nIn my heart\r\nBaby I love you\r\n나 너를 사랑해\r\nAlways love you 곁에 있어 줄래\r\nMy only you 매일 걱정돼요\r\n혹시 사라질까 봐\r\n그댄 나의 사랑 \r\nOnly you', '위로받고 싶은 날 내겐 말해도 돼\r\n위로받고 싶은 날 내겐 기대도 돼\r\n비에 섞인 음악 소리에\r\n혼자 젖은 새벽 감성에\r\n때론 너무 지치고 때론 너무 지쳐서\r\n가끔 아주 가끔\r\n혼술하고 싶은 밤 전화해\r\n혼자 있고 싶은 밤 전 화해\r\n참았던 눈물에 울컥해도 괜찮아\r\n가끔 아주 가끔\r\n혼술하고 싶은 밤 전화 '해\r\n잠이 오지 않는 이 밤에\r\n한 번쯤 소리 내 울컥해도 괜찮아\r\n가끔 아주 가 끔\r\n웃고 싶지 않은 날 그냥 그래도 돼\r\n도망치고 싶은 날 그래 떠나도 돼\r\n비 에 섞인 눈물 소리에\r\n젖어 드는 슬픈 감정에\r\n때론 너무 지치고 때론 너무 지쳐서 \r\n가끔 아주 가끔\r\n혼술하고 싶은 밤 전화해\r\n혼자 있고 싶은 밤 전화해\r\n참 았던 눈물에 울컥해도 괜찮아\r\n가끔 아주 가끔\r\n혼술하고 싶은 밤 전화해\r\n잠이 오지 않는 이 밤에\r\n한 번쯤 소리 내 울컥해도 괜찮아\r\n가끔 아주 가끔\r\n전부 이해할 수 없지만 내게 말해\r\n참았던 눈물이 왈칵 쏟아질까 봐\r\n미치게 답답해 울 컥하는 마음에\r\n가끔 기대도 돼 내게\r\n혼술하고 싶은 밤 외로워\r\n혼자라는 생각

ex. 발라드 Top 100 가사 전체

가사 전처리



#### Kkma

많고 많은 사람 중에 너를 만나서 행복하고 싶어 두 번 다시 울지 않을래 오직 내 눈에는 너만 보여 나를 아껴줘 이제부터 혼자가 아니야 우린 함께니까 나나나 난난 난난나 나나나 난난 난난나

#### Hunnanum

많고 많은 사람 중에 너를 만나서 행복하고 싶어 두 번 다시 울지 않을래 오직 내 눈에는 너만 보여 나를 아껴줘 이제부터 혼자가 아니야 우린 함께니까 나나나 났는 났는다.

#### Okt

많고 많은 사람 중에 너를 만나서 행복하고 싶어 두 번 다시 울지 않을래 오직 내 눈에는 너만 보여 나를 아껴줘 이제부터 혼자가 아니야 우린 함께니까 나나나 났나 났나나 나나나 난난 난난나 나나나 난나 난난나

['그대', '눈', '마음', '순간', '전부', '꿈', '날', '운명', '필요', '내', '결', '2<br/>
<걱정', '나의', '사랑', '미소', '나', '번', '생', '계절', '때', '꽃', '주기', '2<br/>
<너', '위로', '날', '내', '내겐', '건', '기대', '비', '음악', '소리', '혼자', '2<br/>
<젖', '새벽', '감성', '때', '혼', '혼술', '참', '전화', '눈물', '참', '번2<br/>
<', '번쯤', '쯤', '내', '감정', '이해', '수', '마음', '워', '생각', '밤하늘', '2<br/>
<'벌', '너', '내가', '사랑', '내', '결', '어색', '널', '처음', '날', '사람', '숨2<br/>
<', '말', '손', '밤', '하늘', '마음', '표현', '꿈', '웃음', '대도', '겐', '순간2<br/>
<', '영원', '해복', '중', '번', '눈', '나', '이제', '혼자', '니', '난', '당신', 2<br/>
<'나', '워', '대', '마음', '밤', '잠', '매일', '생각', '유치', '심장', '춤', '나', '2<br/>
<'가을밤', '해', '하나', '그대', '오늘', '아름', '그', '고요', '워오', '오', '가을', 2<br/>
<'가음밤', '함', '근간', '내일', '분', '때', '향기', '맘', '지네', '너', '생각', '2<br/>
<'벌', '마음', '순간', '내일', '부담', '날', '산책', '너와', '땐', '보고', '널', 2

ex. 발라드 Top 100 가사 단어

하둡



데이터 크기 ↑

5366	그림
5367	남아
5368	주기
5369	나
5370	년
5371	내
5372	지난
5373	시간
5374	햇살
5375	너
5376	손위
5377	보석
5378	영원
5379	약속

word

id

한글 깨짐 현상











#### 메모장 저장

🧻 ballad.txt - Windows 메모장 파일(F) 편집(E) 서식(O) 보기 그대 눈 마음 순간 전부 1 다가 회색 코 코뿔 뿔 초기 번쯤 쯤 나 눈 이별 그날 모습 거짓말 선명 너 한줄 상태 편 자연 과정 시도 [ 대 곁 여기 두고 그때 모

#### Pom.xml에서 하둡 라이브러리 넣기

```
<dependency>
   <groupId>org.apache.hadoop</groupId>
   <artifactId>hadoop-mapreduce-client-core</artifactId>
   <version>3.1.4
</dependency>
<dependency>
   <groupId>org.apache.hadoop</groupId>
   <artifactId>hadoop-common</artifactId>
   <version>3.1.4
</dependency>
<dependency>
   <groupId>org.apache.hadoop</groupId>
   <artifactId>hadoop-mapreduce-client-jobclient</artifactId>
   <version>3.1.4
</dependency>
```

### 기존 wordcount.java 이용

public static class TokenizerMapper public static class IntSumReducer

[빈도수 2이상 단어 추출] 코드

하둡

가슴속	3	
가요	5	
가운데	2	
가을	3	
가을밤	2	
가지	4	
갈게	3	
감	3	
감정	4	

발라드 Top 100 단어 맵리듀스

가요	2
가운데	2
가지	3
가치	2
간직	2
갈게	3
갈래	4
감	7
감사	2

Ę	낸스
거기	2
거리	11
거울	2
거짓	2
거짓말	3
걱정	6
건	8
걸	3
게	4

어쿠스틱

가요	3
가자	3
가족	2
가지	4
각자	4
간	6
간지	2
갈게	6

길게	0	
힙합		
가까이	2	
가로등	<del>5</del> 3	
가사	4	
가슴	13	
가야	3	
가요	4	
가을	2	
가지	5	
간	2	

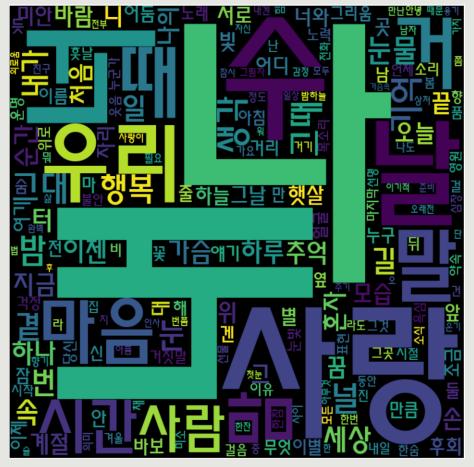
알앤비

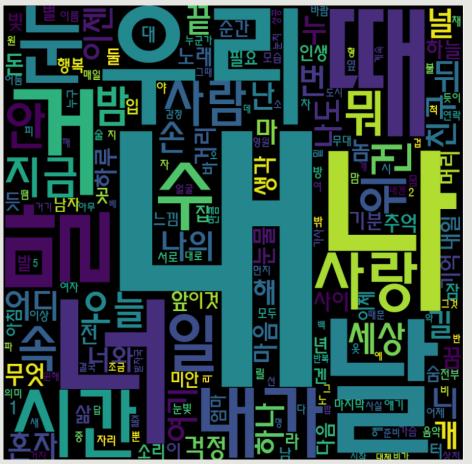
가도 2 가슴 27 가슴속 5 가요 5 가운 2 가지 6 가지마 2 간 4 간직 5

OST

시각화

## Python 워드 클라우드(wordcloud)





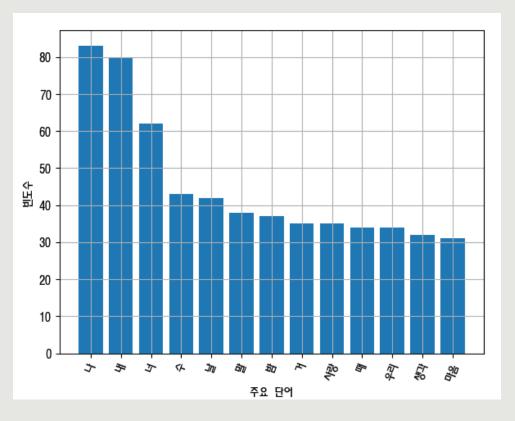
발라드

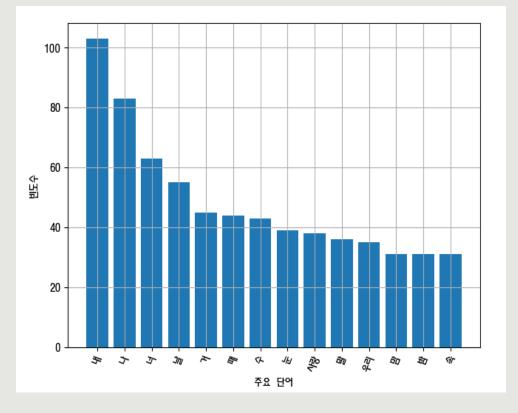
힙합

10

시각화

## ₽ python 그래프(matplotlib)





어쿠스틱

댄스

시각화



여러 개의 문서의 공통된 주제를 찾아내는 프로세스

발라드

댄스

힙합

**OST** 

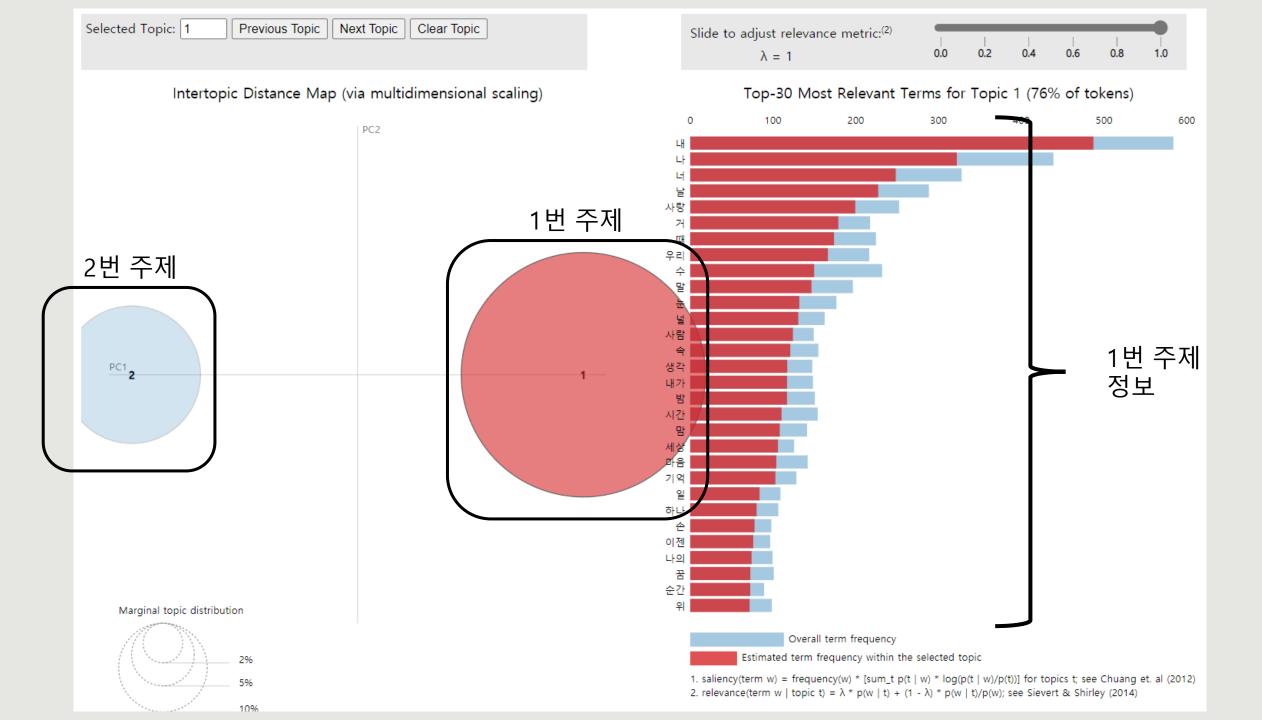
어쿠스틱

알앤비

Input data

- -단어 정수 인코딩
- -단어 빈도수

공통된 주제



## 4. 프로젝트 결론

### 기대효과

- -기존 가사 분석보다 실시간으로 다양하게 분석 가능
- -대중들의 관심과 트렌드를 빠르게 파악 가능
- -나의 음원 목록 분석 등 다른 분야 분석 가능

### 한계점

- -한글 명사 추출이 깔끔하게 되지 않음 (ex. 갈 수 있다 -> '갈', '수' 추출)
- -장르별로 빈도수 높은 단어들이 비슷함