

본 연구에 사용된 모든 코드는 Github<sup>1</sup>에서 **확인 및 결과 재현 가능**

# 자살 유발 정보 디텍팅을 위한 빅데이터 기반 자연어 연구

고귀환 (성균관대학교 인공지능 융합학과 재학생)

E-mail : ie1914@g.skku.edu

Github : <https://github.com/GwiHwan-Go>



# 연구 의도

- 한국 생명 존중 희망 재단에서 자살 유발 정보 모니터링<sup>1</sup>에 관한 활동은 매년 계속해서 이루어지고 있음.
  - 하지만, **완벽한 모니터링은 불가능한 상황**
- 자살 관련 키워드에 대한 정보도 **매우 한정적<sup>2</sup>**
- 자살, 우울에 관한 이야기가 주로 오가는 곳의 정보를 수집해서 중요한 키워드를 찾아내보자.
  - 더 **효과적인 자살/자해 관련 정보 모니터링**

1. 활동에 관한 자세한 정보는 링크 참고 : <https://sims.kfsp.org/usr/main/mainPage.do>

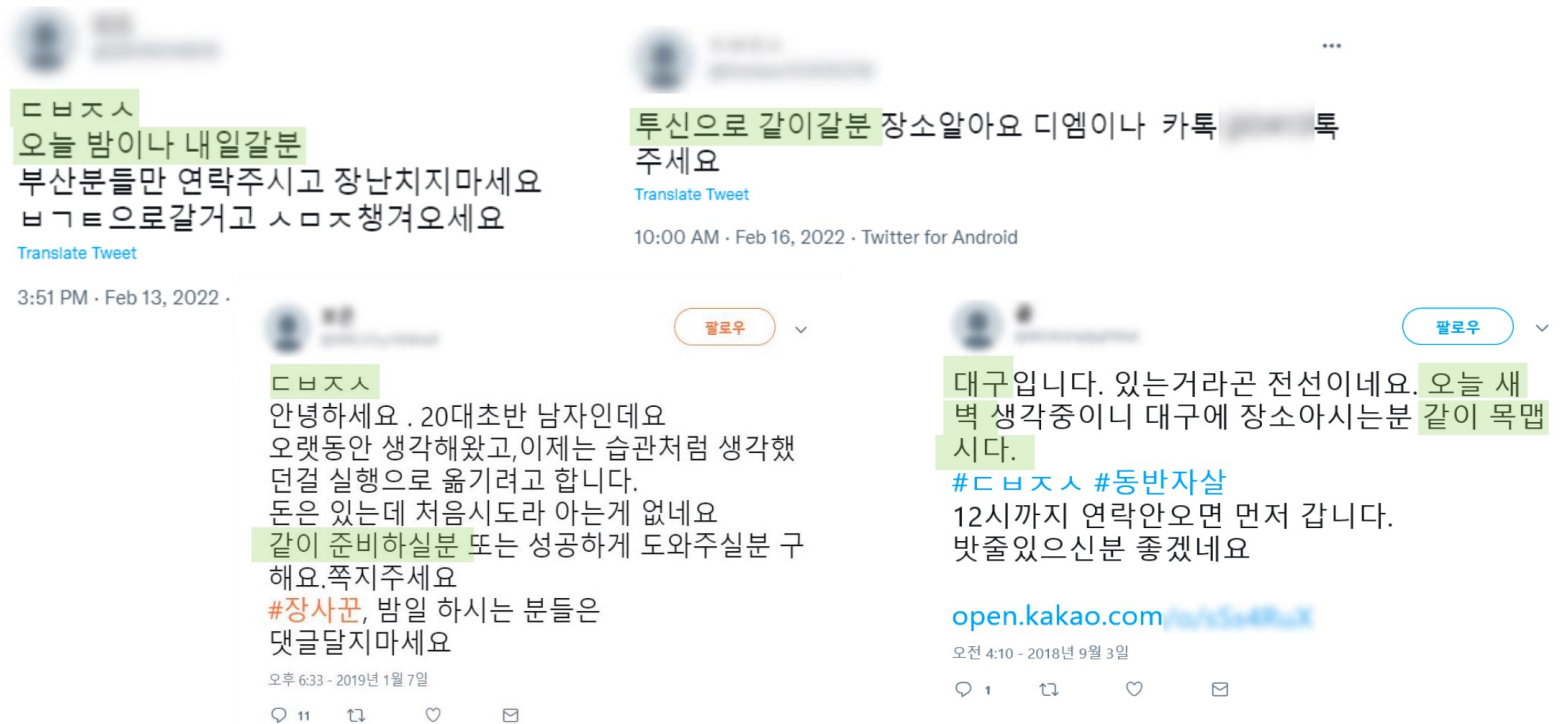
2. 다음 페이지 참고.

# 주로 알려진 자살 관련 정보 키워드<sup>1</sup>

## (2) 검색어 입력

#자살  
#동반자살  
#ㄷㅂㅈㅅ  
#ㅂㄱㅌ  
#자살계  
#죽을래  
#투신  
#자살계  
#자살쇼  
#자살각  
#자해계  
#자해  
#자해전시  
#죽고싶다

⋮



→ **같은 노력에도 제한된 키워드때문에, 일부의 자살 유발 정보만을 모니터링 할 수 있음.**

# 데이터 설명

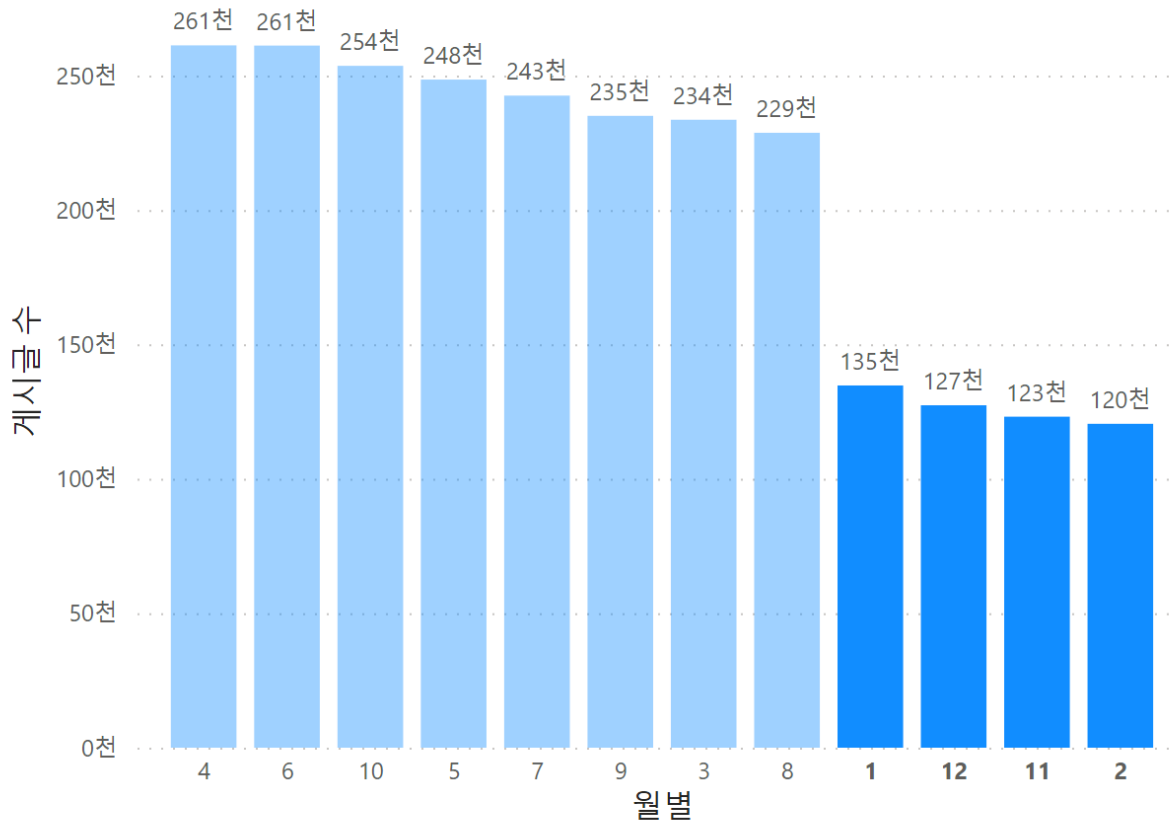
- 사이트 : 디씨 인사이드(dcinside) 우울증 갤러리<sup>1</sup>
- 수집 기간 : 2021.03.03 ~ 2022.10.27
- 선정 이유 : 자살과 관련된 '우울'이란 감정을 가장 잘 내포하고 있을 것이라 추측.
- 수집 방법 : Python 을 이용한 데이터 크롤링<sup>2</sup>
- 수집된 데이터 규모 : 240M

1. [https://gall.dcinside.com/board/lists/?id=depression\\_new1](https://gall.dcinside.com/board/lists/?id=depression_new1)

2. 코드는 깃헙 리파지터리에서 확인가능([https://github.com/GwiHwan-Go/Detect\\_SC](https://github.com/GwiHwan-Go/Detect_SC))

# 월 별 데이터 분포 시각화

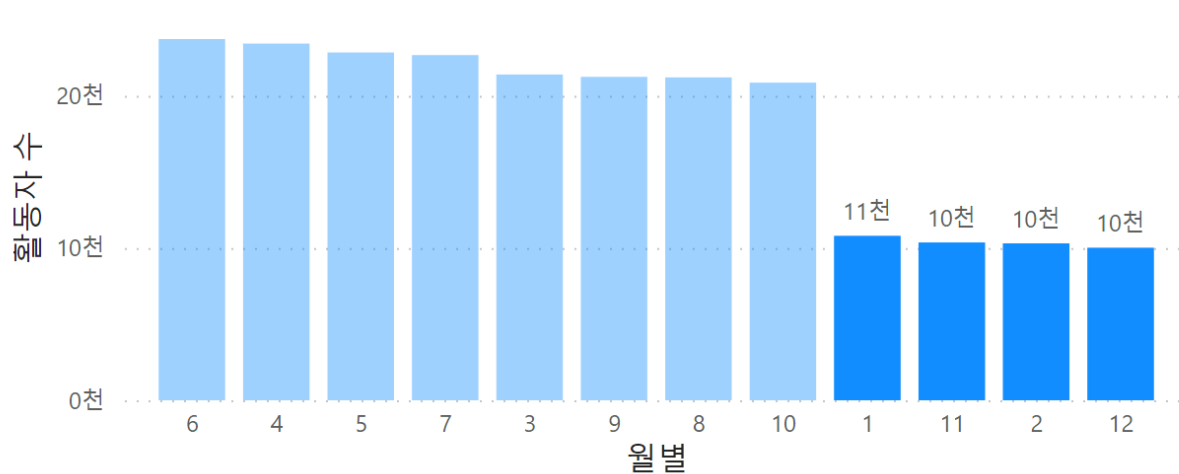
게시글 수 및 합계 date개, 월별



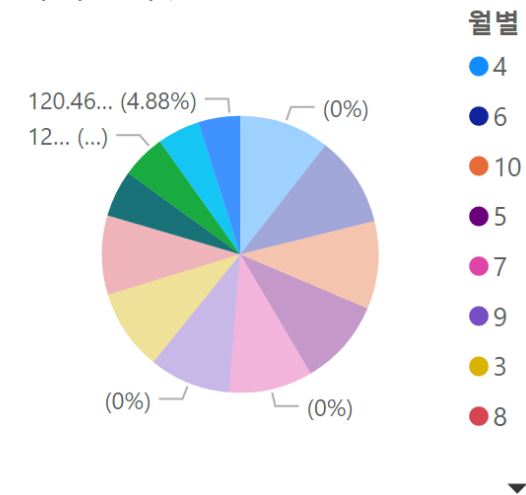
수집된 데이터의 기간 : **21.3.3 ~ 22.10.27**

11,12,1,2 월은 21년도와 22년도 사이에 겹치는 범위가 아니다.  
이를 보여주듯, 11,12,1,2월의 데이터는 다른 달의 데이터에 비해 약,  
50% 적다.

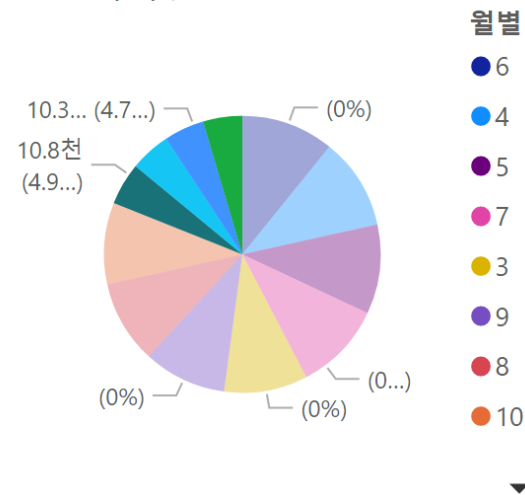
활동자 수, 월별



게시글 수, 월별

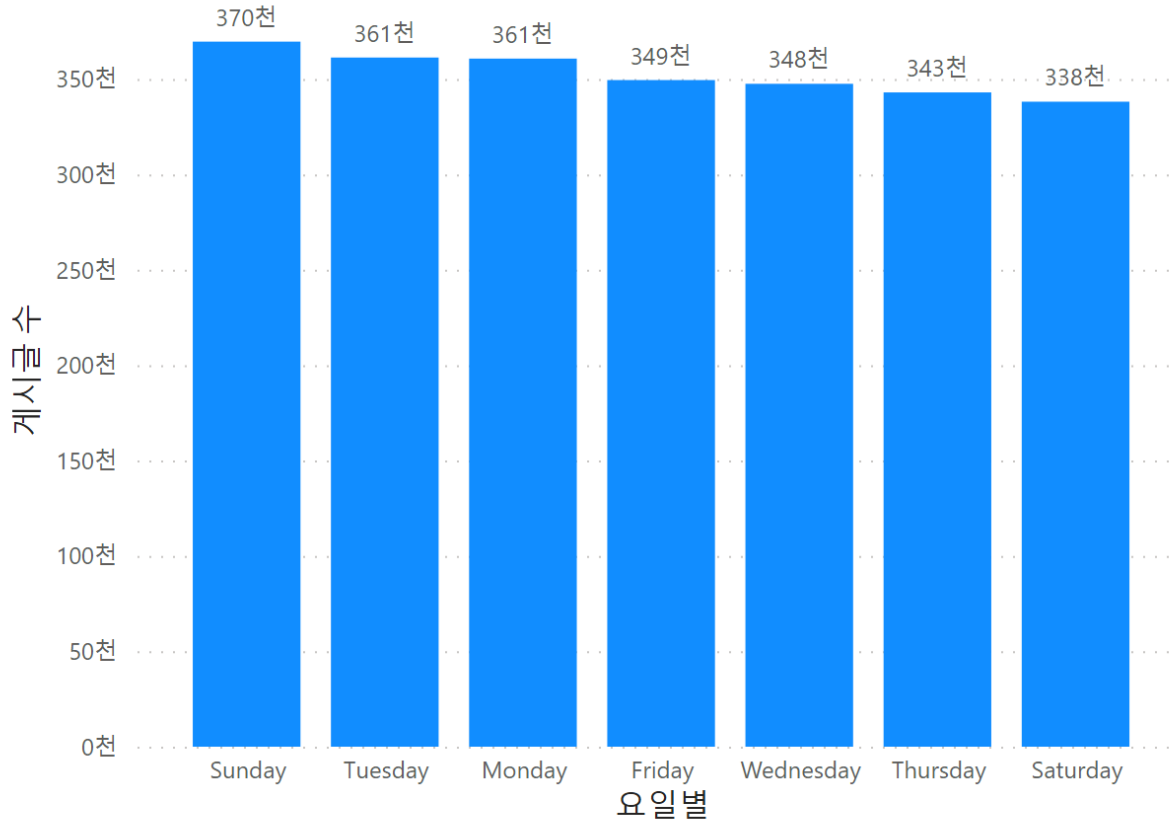


활동자 수, 월별



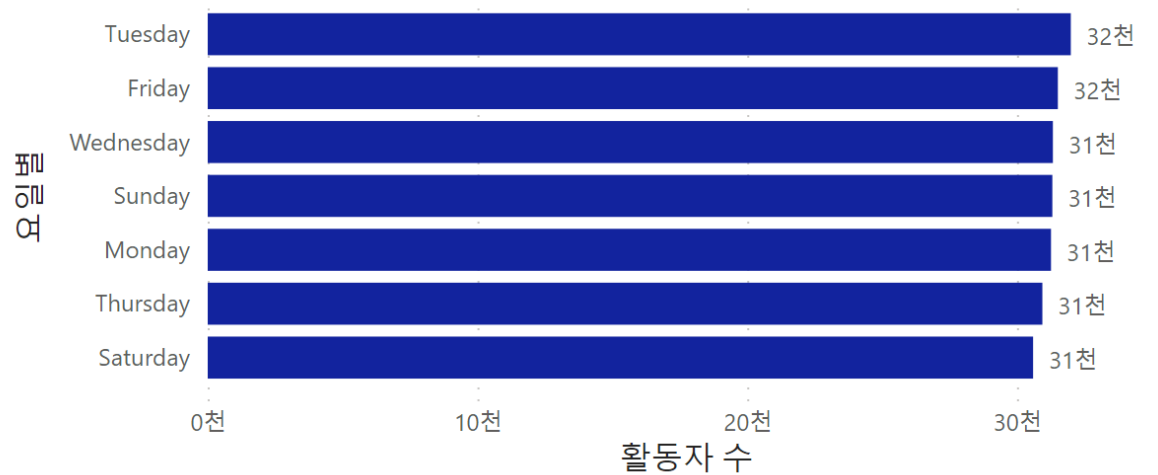
# 요일별 데이터 분포

게시글 수, 요일별

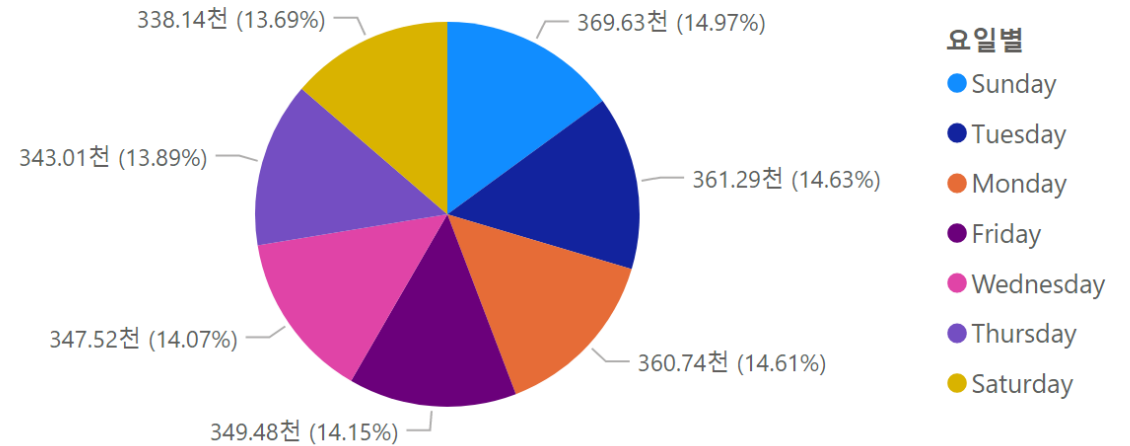


일요일이 다른 요일보다 게시글의 수가 더 많음을 알 수 있음.  
하지만, [요일 별 활동자 수] 그래프를 보면, 일요일의 활동자 수는 많은 편에 속하지 않음.  
==> 일요일엔 특정 유저들의 활동이 활발함을 알 수 있음.

활동자 수, 요일별



게시글 수, 요일별



# 데이터 분석 전문 용어 정리

- 단어 정제(preprocessing) : 조사, 숫자, 특수문자 제거
- 토큰화(tokenizing) : 문장을 형태소 혹은 단어로 나누는 행위, 여기서 나뉜 형태소 혹은 단어를 토큰(token)이라 함.
- 신조어 : 사전에 등재되지 않은 단어
- 비지도 학습방식 : 데이터 속에서 컴퓨터가 스스로 규칙을 찾아 학습하게 하는 방법
- 워드 클라우드 : 단어를 빈도수 별로 강조한 시각화 도구

# 자살 관련 키워드 추출 분석 방법<sup>1</sup>

- 텍스트 전처리

- 신조어의 추출을 위해 사전 기반이 아닌 SoyNLP의 비지도 학습방식의 토큰화 도구 이용해서 수집한 데이터를 학습시킴.
- 실제 사용 예 : "야이개색야" – 기존 토큰라이저 : '야', '이개', '색', '야' | 본 연구 : '야', '이', '개색', '야'

- 텍스트 분석 방법

- Positive Point Mutual Information(PMI) 이용
- 단어와 문맥의 상관성을 측정하는 방법

주의 : 본 연구 결과에는 비속어/저속어가 다수 포함되어 있습니다.

자료 출처 : [https://github.com/lovit/soynlp/blob/master/tutorials/pmi\\_usage.ipynb](https://github.com/lovit/soynlp/blob/master/tutorials/pmi_usage.ipynb)



**-정제된(preprocessed) 단어들이임.**

## 전체 데이터에 관한 워드클라우드

사람	새끼	존나	jpg	씨발	여자	하는	갤러	너무					
..	근데	??	친구	ㅌㅌ	아니	시발	나도	ㅍㅍㅍ	ㅈㄴ	남자			
ㅋㅋ	생각	보고	그냥	요즘	이유	노래	는데	먹고	담배	얼굴	..	ㅎ...	
진짜	인생	시간	들어										
	사랑	기분	추천										
	하면	보면	인스타	언제									
○○	지금	같음	있음	카톡	한번								
	병신	싫다	연애	들이	아침	자꾸							
		우울증	해서	님들	어제	하...							
...	같은	ㅁㅈㅈ	....	다시	소리	서울	이...						
ㅅㅅ	하고	이제	많이	남친	홍대	학교	하실						
	사진	여기	누구	여친	해야	한남	보다						
	머리	좋아	하루	좋은	우울	있냐	운동	30					
오늘	ㄹㅇ	우리	있는	누나	없음	ㅋ...	없어						

# 연구 결과 설명

- 분석의 시작이 되는 단어는 3페이지의 단어들로 구성함.
- 문맥 단어
  - 특정 단어와 같이 등장한 단어.
  - => 같이 쓰이는 단어들이 어떤 것이 있는지
- 유사 단어
  - 같은 문맥을 가지기 때문에 유사하다고 여겨지는 단어.
  - => 어떤 단어들이 대체해서 쓰이는 지, 같은 맥락에서 쓰이는지
- 시각화 된 결과를 곱씹어 볼 수록 많은 것을 느낄 수 있을 것이다.

주의 : 본 연구 결과에는 비속어/ 저속어가 다수 포함되어 있습니다.



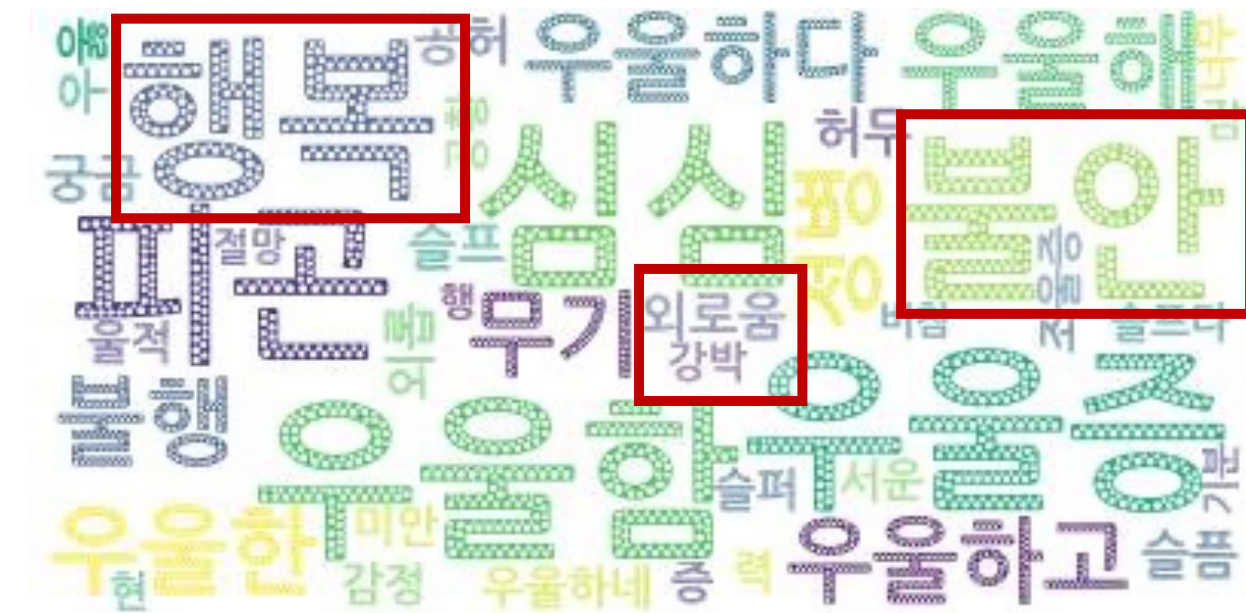
동반은 주로 '입대'라는 단어와 연관을 보이며, '자살'과는 큰상관을 보이지 않는다.

**'자살'이란 단어는 '살고싶다'라는 단어와 연관이 크다  
사실 그 누구보다 살고 싶어하는 사람들인 것이다.**





우울: 문맥 단어↑ 유사 단어↓



살자: 문맥 단어↑ 유사 단어↓

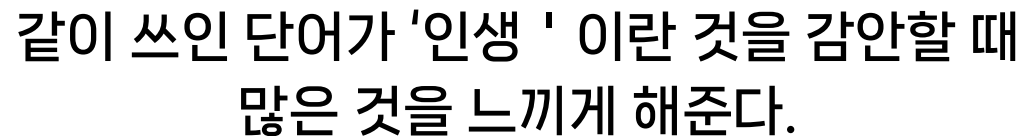
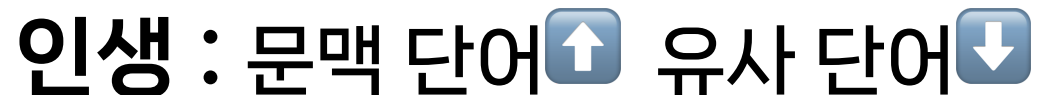






따라서 커뮤니티 회원들은 이를 변형 시켜서 사용한다.

본 분석은 커뮤니티 회원들이 변형 시켜서 사용하는 단어들을 성공적으로 추적하였다. [죯해, 자혜, ㄱㅎ, 좇해]





## 죽고싶다: 문맥 단어 유사 단어

굶어아안 올 어머 강 에 나 정말  
 니미 번만 사귄 귀찮은데 안아 생겨서 그냥 ??? 차라리  
 죽을래 나가 노세인 죽어  
 없으면 사겨 없을 우울해서 안 갇 으아 버릴래  
 짹 안 아무래도 싫 채 라멘

**죽을래** : 문맥 단어  유사 단어 

버릴래  
필라다

## 죽고싶다의 다른 표현으로

'버릴래', '죽구싶다' 도 사용되고 있다.

**비표준어**는 모니터링에 가장 방해가 되는 요소 중 하나이다.

본 분석은 비표준어 문제를 잘 해결하고 있다.





# 워드 클라우드 결과 종합

- '자살, 죽음, 자해' 등이 '편안, 행복, 살고 싶다.' 등의 긍정적 단어와 같이 등장한다.
  - 모니터링에서 '살고 싶다', '편안', '행복' 등의 키워드도 고려해볼 필요가 있다.
  - 죽을 생각이 없는 사람은 '살고 싶다.' 라는 생각을 하지 않기 때문이다.
- 자동 검열을 회피하기 위해서 사용되는 비표준 단어 추적
  - ⇒ 자해 대신 [쫓해, 자혜], 자살 대신 '자ㅏㅏㅏㅏ', 죽고싶다 대신 '죽구싶다' 등의 의도적인 비표준어 사용을 성공적으로 추적.
  - ⇒ 연구결과에 등장한 상식적으로 이해가 되지 않는 단어들 역시 특정 단어의 변형 형태일 가능성이 있다.



# Graph 데이터 설명

- Graph 생성 알고리즘

```
until recursion boundary :  
  Seed -> Context_word_list  
  for Context_word in Context_word_list  
    build graph
```

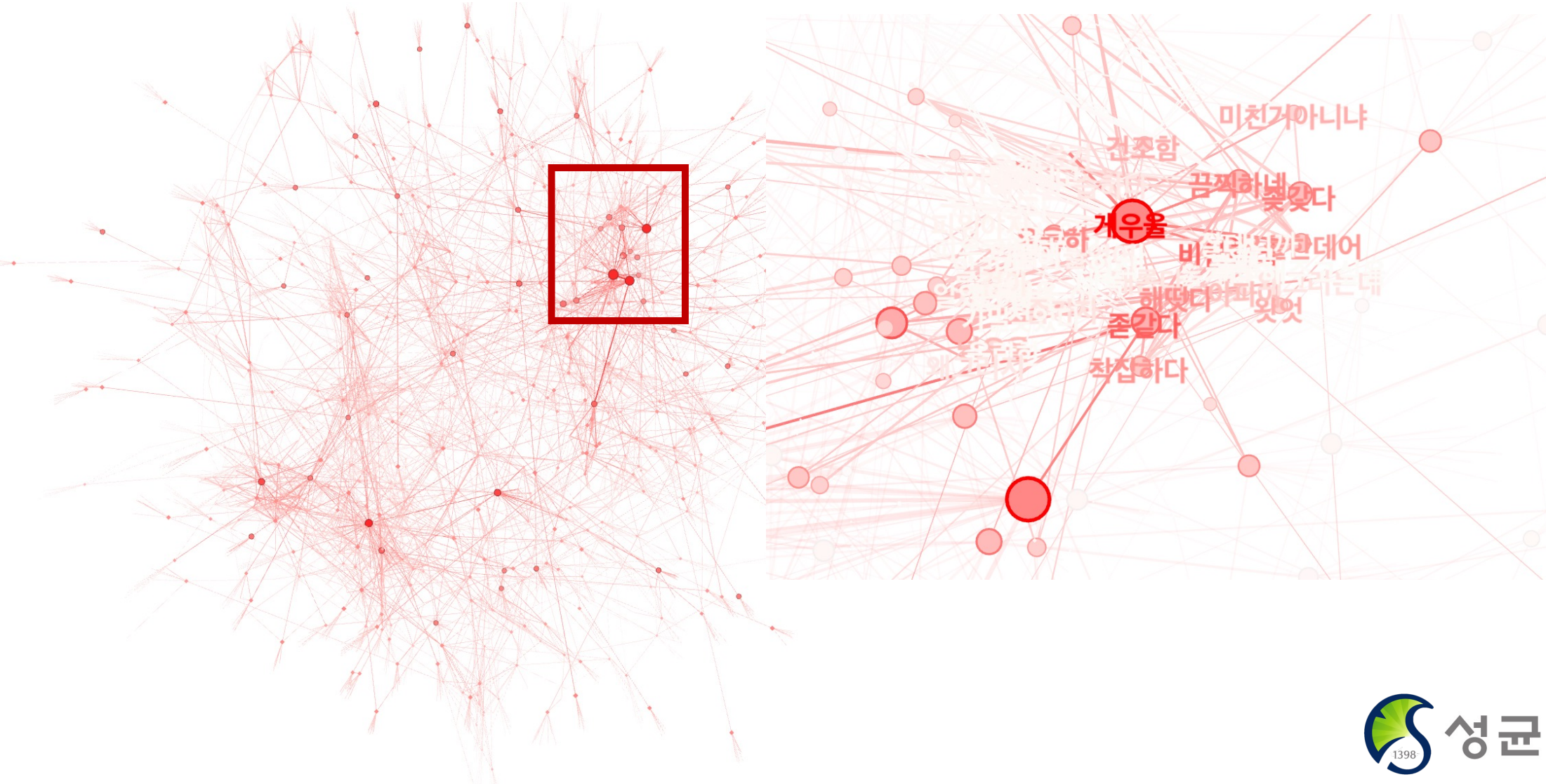
- 그래프 전문 용어

- 페이지 랭크<sup>1</sup> : 중요도를 표시하는 방식
- 노드 : 단어를 가리킴
- 엣지 : 단어와 단어 간의 연결선

# 그래프 데이터 설명

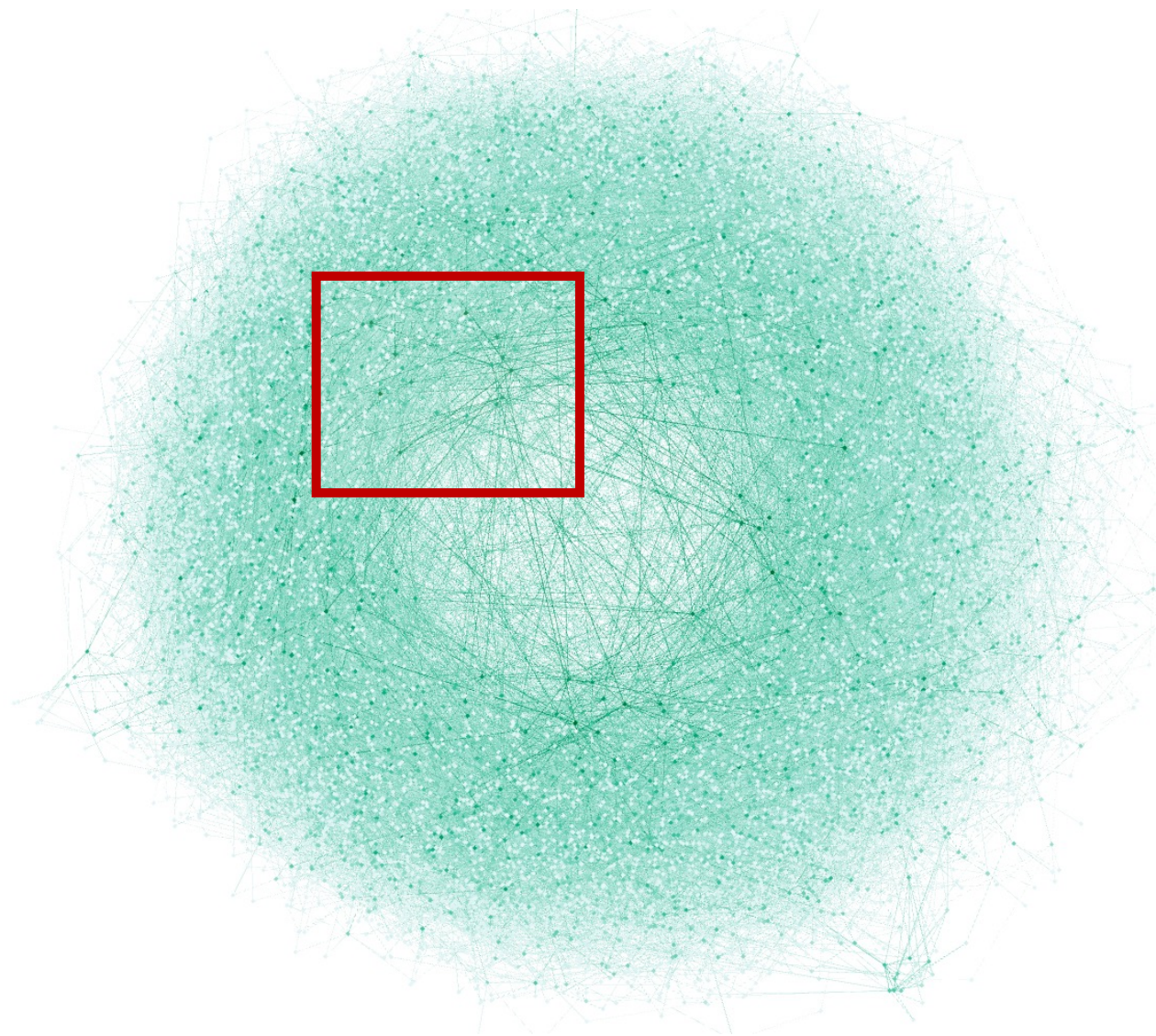
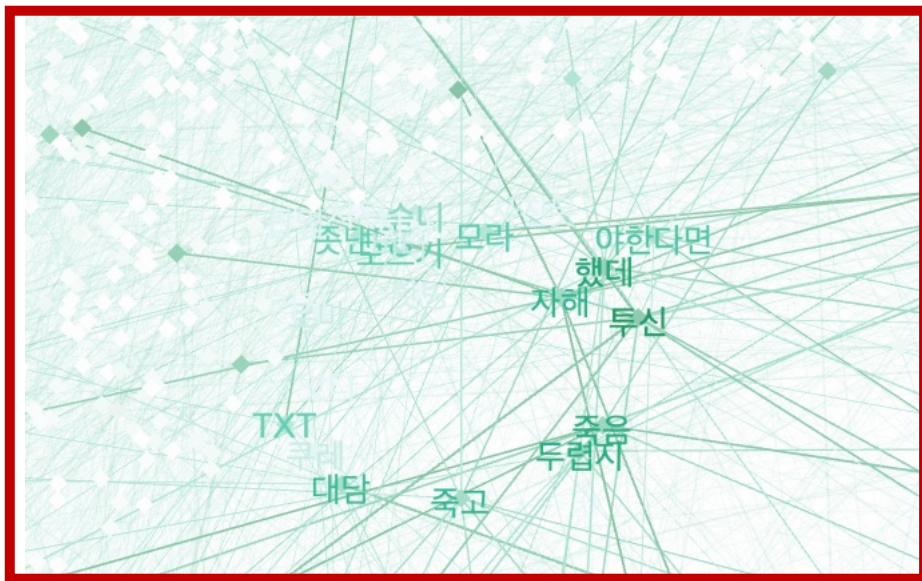
- 문맥 단어 그래프
  - 단어와 문맥 단어들 간을 연결한 그래프
- 유사 단어 그래프
  - 단어와 유사 단어들 간을 연결한 그래프
- 각, 그래프당 4만 여개의 노드와 20만 여개의 엣지가 존재.
- 효과적 분석을 위해서 페이지 랭크를 기준으로 필터링 수행.

# 데이터 오버뷰 (유사 단어 그래프)





# 데이터 오버뷰 (문맥 단어 그래프)



# 그래프 데이터 활용 방안

1. 유사 단어 그래프와 문맥 단어 그래프를 조합
2. 아래의 테스트를 수행하는 분류기 빌드
  - 특정 게시물, 혹은 작성자의 우울 지수 평가
  - 자살 관련 정보 분류기



# Unique Origin Unique Future



자세한 코드는 깃헙에서 확인 및 연구결과 재현 가능  
[https://github.com/GwiHwan-Go/Detect\\_SC](https://github.com/GwiHwan-Go/Detect_SC)