

KU BIG DATA

# TEXT MINING

최종 프로젝트 발표

김강우 배건희 문선미 박기찬 이세희

# KU BIG DATA TEXTMINING

---

## **01** Study review

## **02** Prepre-processing

- Crawling
- 각종 처리

## **03** Pre-processing

- Stopwords 제거
- 토큰 규모 축소
- 문서 규모 축소
- Exception

## **04** Modeling analysis

- LDA
  - Word cloud
-

TEXT MINING

# STUDY REVIEW

- 주요 개념

## [ 숨어있는 정보를 찾아내는 Text Mining ]

The image shows a screenshot of the Twitter profile for the UN Spokesperson (@UN\_Spokesperson). The profile header features a blue background with the text "Data hidden in plain sight" and the United Nations logo. The profile information includes the name "UN Spokesperson", the handle "@UN\_Spokesperson", the bio "Official Twitter account of the Office of the Spokesperson for United Nations Secretary-General Ban Ki-moon.", the location "New York, USA", the website "un.org/sg/spokesperso...", and the join date "Joined May 2010". The profile also shows a "Tweet to UN Spokesperson" button and a section for "3,008 Photos and videos". The tweets section displays three tweets. The first tweet is about maintaining unity in tackling security challenges on the Korean Peninsula. The second tweet is about ethics being built into the ideals and objectives of the United Nations. The third tweet is about a ban on Ambassador Joseph V. Reed. Annotations with green boxes and arrows point to various elements: "Author" points to the profile name, "Description" points to the bio, "Location" points to the location field, "Tweet" points to the tweet text, and "Time" points to the timestamp.

**Author**

**Description**

**Location**

**Tweet**

- Topic
- Sentiment

**Time**

## [Normalization]

“ 대문자를 모두 소문자로

“ ‘ ’ 를 기준으로 split한 결과

## [Stemming]

“ Word의 root word, root from을 찾는 것(어간 추출)

“ 단어 그 자체만을 고려

“ 코딩 결과도 words 1에 있는 단어들의 root word인 list를 추출

**Normalization  
And  
Stemming**

## [ Text Classification ]

“Input text에 대해서 Labeling 하기

Topic Identification : Politics, Sports, Technology?

Spam detection : Spam or not?

Sentimental analysis(감정분석): Movie review -> positive or negative?

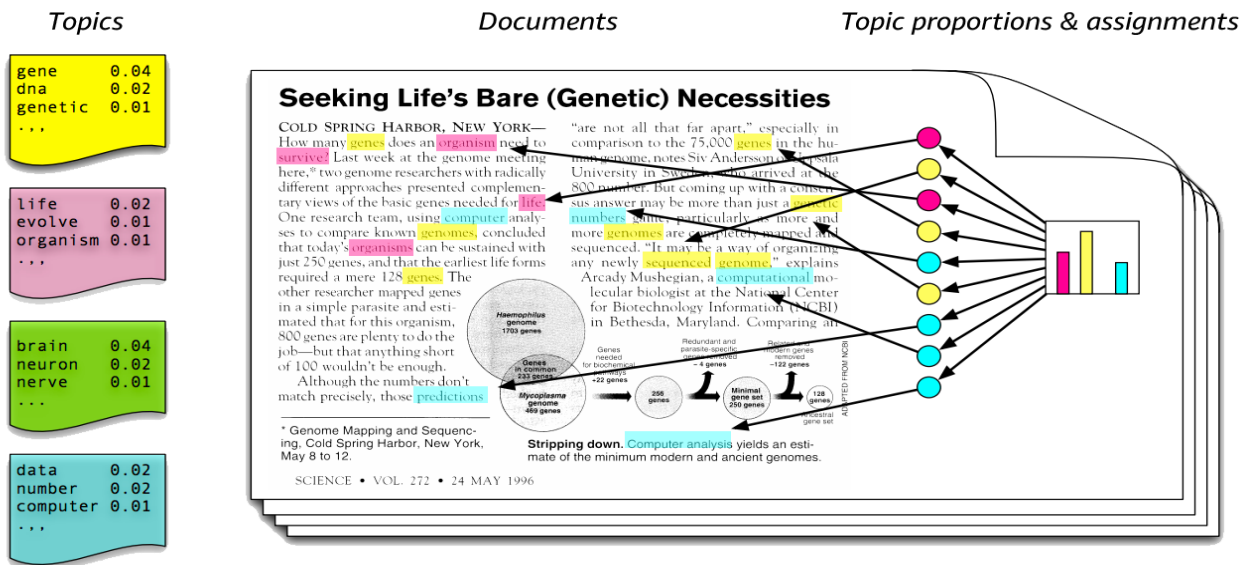
Spelling Correction : Weather or Whether? Color or Colour?

-> Text Classification도 여태까지 알아왔던 Classification과 마찬가지로  
'Supervised Classification' !

## [ Bag of Words ]

- “ Document를 자동으로 분류하기 위한 기법  
글에 포함된 단어(Word)들의 분포 -> 이 문서가 어떤 종류의 문서인지를 판단
- “ Bag-of-words : Term- Document Matrix  
어떠한 text에 포함되어 있는 words에 대해, 순서에 상관없이, 각 words의 사용 빈도  
값으로, 바꾸어 주는 방법.
- “ Bag-of-words : Representation in a Vector Space  
궁극적으로 Frequency of words를 보고, 그 Text의 Contents를 추론하는 방법  
많이 사용된 Words가 그 Text의 핵심 단어일 가능성이 높으므로.

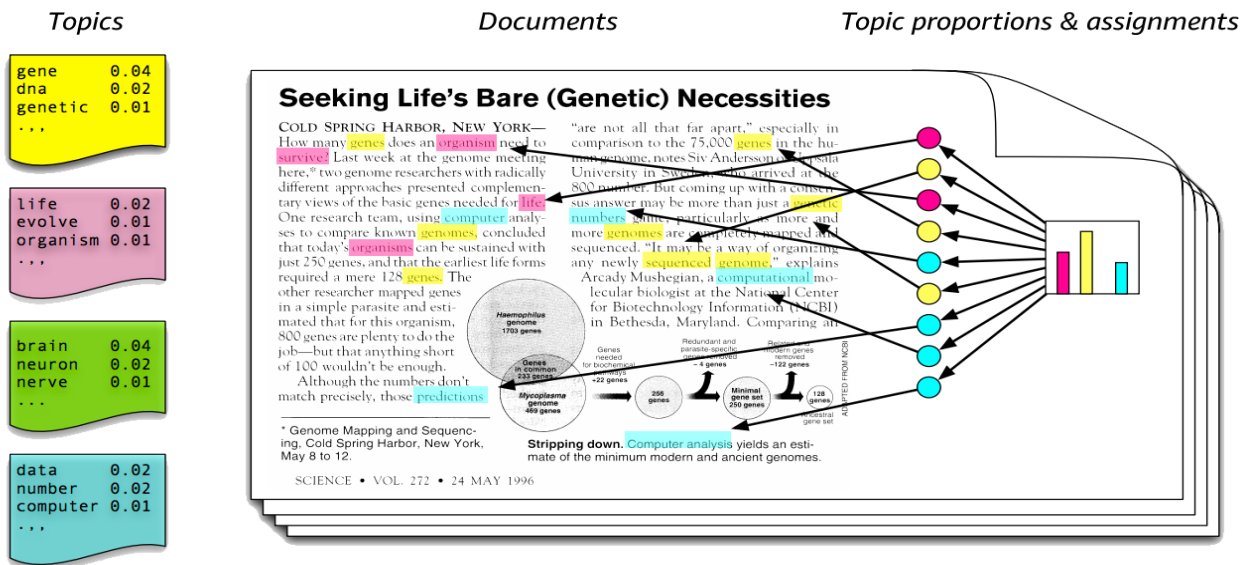
# [ LDA (Latent Dirichlet Allocation) ]



“ 문서와 토픽, 단어와 토픽 간의 잠재(Latent)적인 관계들을 디리클레 분포(Dirichlet Dist.)로 나타내고 이들을 활용해 토픽 할당(Allocation)



# [ LDA (Latent Dirichlet Allocation) ]

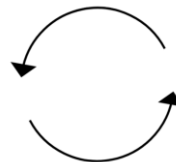


“ 문서와 토픽, 단어와 토픽 간의 잠재(Latent)적인 관계들을 디리클레 분포(Dirichlet Dist.)로 나타내고 이들을 활용해 토픽 할당(Allocation)

## [ LDA (Latent Dirichlet Allocation) ]

### 알고리즘

1. 각 단어에 Random하게 Topic 배정 후 분포 생성
2. 문서와 토픽 / 단어와 토픽 간 디리클레 분포를 추정
3. 단어 별 결합분포를 활용해 단어별 Topic 재배정



수렴할 때까지 **반복**

## [ Named Entity Recognition task ]

어디까지 객체로 인식할 것인지 사전 결정

The patient is a 63-year-old female with a three-year history of bilateral hand numbness and occasional weakness.

Within the past year, these symptoms have progressively gotten worse, to encompass also her feet.

bilateral hand numbness

bilateral hand numbness

hand



TEXT MINING

# PREPREPROCESSING

- Crawling
- 각종 처리

## Crawling

 연세대학교 대나무숲

5월 1일 오후 9:18 · 🌐

연대숲 #59299번째 외침:

2018. 5. 2 오전 10:33:16

대—하

대숲!!!!!!!!!!!!!!!!!!!!

남자친구가 인스타에서 다른 여자 인스타스타들을 팔로우하는게 싫어  
요!!!!!!!!!!!! 하나둘도 아니고!!!!!!!!!!!!

물론 예뻐요 제가 봐도!!!!!!

근데 싫어요 그냥 다 싫어!!!!!!!!!!!!

나만!!!!!!!!!!!! 싫은건가!!!!!!!!!!!!

남자친구야 이거 꼭 봐주길바라 그냥 그렇다구<(˙˘˙)>

휴 고마워요 대숲ㅎㅎㅎ



연대숲 #59299번째 외침:

2018. 5. 2 오전 10:33:16

대—하

대숲!!!!!!!!!!!!!!!!!!!!

남자친구가 인스타에서 다른 여자 인스타스타들을 팔로우하는게  
싫어요!!!!!!!!!!!! 하나둘도 아니고!!!!!!!!!!!!

물론 예뻐요 제가 봐도!!!!!!

근데 싫어요 그냥 다 싫어!!!!!!!!!!!!

나만!!!!!!!!!!!! 싫은건가!!!!!!!!!!!!

남자친구야 이거 꼭 봐주길바라 그냥 그렇다구<U+1559>(U+  
2022><U+0300><U+2038><U+2022><U+0301><U+2036>  
<U+1557>

휴 고마워요 대숲ㅎㅎㅎ

// Rfacebook 라이브러리를 이용 .csv로 크  
롤링

```
yonsei['message'] = yonsei['message'].str.replace(r'<ed>(<U[A-Z0-9]+>)+','')
```

## 기본문장처리

```
대—하
대쉴!!!!!!!!!!!!!!!!!!!!
남자친구가 인스타에서 다른 여자 인스타스타들을 팔로우하는게
싫어요!!!!!!!!!!!! 하나둘도 아니고!!!!!!!!!!!!
물론 예뻐요 제가 봐도!!!!!!
근데 싫어요 그냥 다 싫어!!!!!!!!!!!!
나만!!!!!!!!!!!! 싫은건가!!!!!!!!!!!!

남자친구야 이거 꼭 봐주길바라 그냥 그렇다구

휴 고마워요 대쉴ㅎㅎㅎ
```

### // 맨앞 공백 제거

```
yonsei['message'] =
yonsei['message'].str.lstrip()
```

### // 문장을 마치는 부호들을 !로 바꿔주기 (.은 다루기 어려워서 !로 대체)

```
yonsei['message'] = yonsei['message'].str.replace(r'[.~]', '!')
yonsei['message'] = yonsei['message'].str.replace(r'[.]', '!')
yonsei['message'] = yonsei['message'].str.replace(r'[?]', '!')
yonsei['message'] = yonsei['message'].str.replace(r'[~]', '!')
```

```
yonseil['message'] = yonseil['message'].str.replace(r'ws+', '!')
yonseil['message'] = yonseil['message'].str.replace(r'!', '! ')
```



## 세부처리

대—하! 대숲! 남자친구가 인스타에서  
다른 여자 인스타스타들을 팔로우하는게  
싫어요! 하나둘도 아니고! 물론 예뻐요  
제가 봐도! ! 근데 싫어요 그냥 다 싫어!  
나만! 싫은건가! 남자친구야 이거 꼭 봐  
주길바라 그냥 그렇다구 휴 고마워요 대숲  
ㅎㅎㅎ

### // 기타 문장부호, 초성 제거

```
yonsei['message'] = yonsei['message'].str.replace(r"^[가-힝0-9a-zA-ZWs.,/!~]",
")
```

### // 대—하! 대숲!같은 의미없는 인사말 제거

```
yonsei['message'] = yonsei['message'].str.replace(r'^(\대—[가-힝!])' ,")
yonsei['message'] = yonsei['message'].str.replace(r'^(\ 대[가-힝!])' ,")
```

토큰화: R의 **extractNoun**을 이용하여 명사화하여 토큰으로 만든다.

#### MESSEGE

남자친구가 인스타에서 다른 여자 인  
스타스타들을 팔로우하는게 싫어요!  
하나둘도 아니고! 물론 예뻐요 제가  
봐도!! 근데 싫어요 그냥 다 싫어!  
나만! 싫은건가! 남자친구야 이거  
꼭 봐주길바라 그냥 그렇다구 휴 고  
마워요 대썬



#### TOKEN

['남자친구', '인스타', '여자', '인스  
타스타들', '팔로우하는', '하나둘', '  
저', '근데', '나', '것', '남자친구', '  
이거', '봐주길바']

TEXT MINING

# PREPROCESSING

- STOPWORDS 제거
- 토큰 규모 축소
- 문서 규모 축소
- Exception

## STOPWORDS 제거

### // STOPWORD?

인터넷 검색 시 검색 용어로 사용하지 않는 단어. 관사, 전치사, 조사, 접속사 등 검색 색인 단어로 의미가 없는 단어이다.

#### 예시

여러분, 저, 것, 다들, 전, 나, 그거, 애, 기,  
너, 내, 누군가, 진짜, 데, 제, 분, 거, 수,  
의, 은, 는, 이, 가, 들, 중, 듯, 네게, 재, 지,  
사, 줄, 절, 로, 안, 못, 몇번, 너희, 이젠,  
때문, 도, 터, 이후, 관리자님, 안녕, 명,  
애, 지금, 우리, 언제, 이번 등

한국어 불용어 리스트

+

대나무숲 특성에 맞게 추가한 불용어

```
# with stop_lists
texts = [[word for word in doc if word not in stop_lists] for doc in documents]
```

## 빈도수 낮은 단어 제거

```
from collections import defaultdict
frequency = defaultdict(int)
for text in texts:
    for token in text:
        frequency[token] += 1

texts = [[token for token in text if frequency[token] > 10]
         for text in texts]
```

// 10번 이상 등장한 단어만 남기기

## 토큰수 적은 문서 제거

```
texts = [doc for doc in texts if len(doc) > 5]
```

#21591번째사자후  
16학년 여러분 23살 누나 어때요?ㅋㅋㅋㅋ 너무 늙었나...

#21590번째사자후  
아무리 관심없다 관심없다해도 자꾸 생각나고 보고 싶은걸 어떡해...

// 토큰 5개 이상으로 구성된  
문서만 남기기

## [Exception]

### “한글파괴 빌런

#21489번째사자후

역석 한양대다. 지금 보시는 아시게찌마는 한양인의 시험준비 메카니즘은 상당히 조크드요? 보세요. 긴장없이 푸는데 에이뽀를 두르냈스니까. 즈응말 공부 잘해요. 제가 그딴 말씀드리지마는 즈른 불리한 상황에스 보닌 샤뿌를 휘두를 수 있는 손슈가 그롭게 만치 안타. 그런데 보닌 소느로 학점을 맨드러 내쓰니까 그즈을 칭찬하지 아눌 수가 없을고 가타요. 즈는 이롭게 평가를 해요. 췌게 최종상급 손슈다. 그롭기 때문에 인뿌라가 중요흔그에요. 그롭기 때문에 돔수업실이 이쓰야 해요. Thumbs up

### “2 빌런

#22222번째사자후

17학번2 들어오니 새내기적 2야기가 생각나네요..ㅋㅋㅋ

새내기 때 페2스북을 안해서 대나무숲을 몰랐을 적

저는 학교안에 대나무숲2 정말로 2씨는 줄 알았습니다!

자신의 고민을 포스트2스에 적어서 대나무에 붙2면

사람2 없는 저녁에 대나무숲지기들2 슬금슬금 와서 포스트2스를 ㅁ

ㅋㅋㅋ

한양2니들은 그 주제에 대해서 서로 얘기도 나누고요!

어떤사람들은 고민을 붙2다가 눈2 맞아서 연애를 하고(제 소원..)

선봉고발을 하다가 머리채를 잡는 줄 알았어요..ㅋㅋㅋ



삭제

## [ Exception ]

“ ?

#32579번째올림

2015. 06. 16. 오후 02:09:00

< 성균관대학교 대나무숲 >

[illegible]

“ 띄어쓰기 실종

[전공진입후 달린다다꺼 저내가갈거]



# 삭제

## “ 지나친 감정표현

#20442번째사자후

내일부터 실습시작이에요!!!까르르르르카하하하하하함함함  
거의 3달간 폭 쉬고 오랜만에 하려니까 너무너무 가기 두렵네여  
우리와 4학년 모두 화이팅함함해해함함함  
#케이스 #퀴즈 #술기 #모스비

TEXT MINING

# MODELING ANALYSIS

– LDA



## [ 학교별 포스트 수 ]

```
In [87]: number_dic
```

```
Out [87]: {'center': 3096,  
          'hanyang': 3129,  
          'konkuk': 2087,  
          'korea': 3304,  
          'kyunghee': 3167,  
          'seoul': 4764,  
          'sirip': 2914,  
          'skk': 3345,  
          'sugang': 3182,  
          'yonsei': 3147}
```

## [ 기존 Description ]

```
In [67]: len(dictionary.token2id)
```

```
Out [67]: 161125
```

```
In [86]: len(texts2)
```

```
Out [86]: 32135
```

```
In [23]: dictionary.num_pos
```

```
Out [23]: 1717748
```

## “ # of Unique Token

```
In [235]: len(dictionary.token2id)
```

```
Out [235]: 9412
```

## “ # of Total Posts (row)

```
In [236]: dictionary.num_docs
```

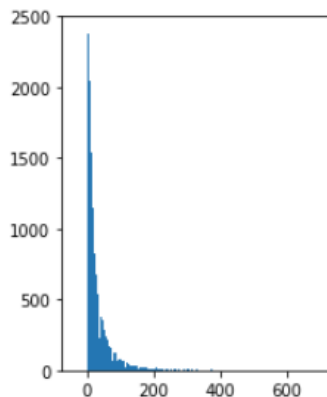
```
Out [236]: 25738
```

## “ # of Total Tokens

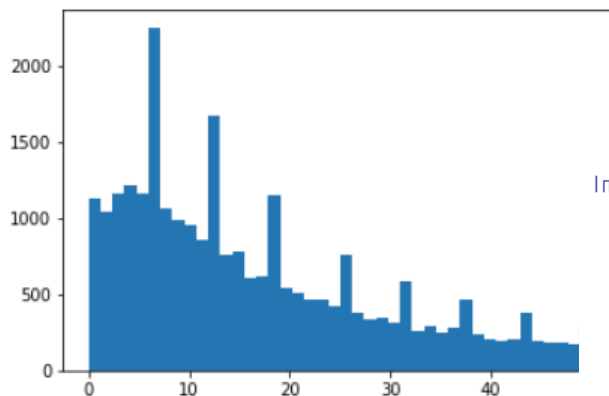
```
In [237]: dictionary.num_pos
```

```
Out [237]: 926428
```

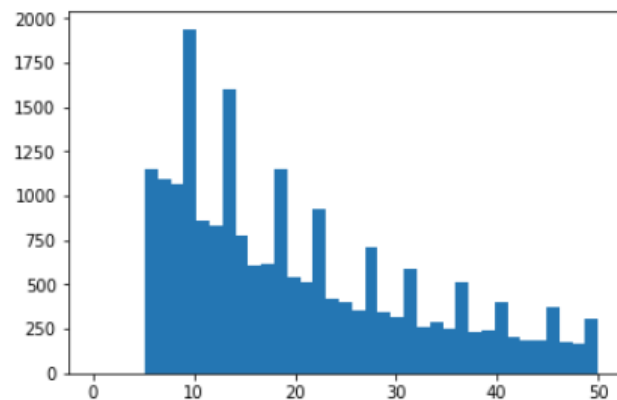
## [ Token 수에 따른 Post 수 Histogram ]



```
In [106]: plt.hist(len_doc, bins = 'auto', range = (0, 50))
plt.show()
```



```
In [109]: plt.hist(len_doc, bins = 'auto', range = (0, 50))
plt.show()
```



전처리 후 ->

## [ K(Cluster 수) 에 따른 Perplexity 값 ]

```
In [239]: grid_k = [3, 5, 7, 10, 15, 25, 30, 40]
          vec_perplex = []
```

```
In [202]: for i in np.arange(len(grid_k)):
          lda_model = gensim.models.LdaModel(corpus, num_topics=grid_k[i], id2word = c
          perplex = lda_model.log_perplexity(corpus)
          vec_perplex = vec_perplex + [perplex]

          print(i, "th iteration is done")
          print("perplexity is ", perplex)
```

```
0 th iteration is done
perplexity is -7.43108611888
1 th iteration is done
perplexity is -7.39757952516
2 th iteration is done
perplexity is -7.43762033601
3 th iteration is done
perplexity is -7.43255894301
4 th iteration is done
perplexity is -7.45502740344
5 th iteration is done
perplexity is -7.51579825474
6 th iteration is done
perplexity is -7.53918962392
7 th iteration is done
perplexity is -7.60939018584
```

$$Perplexity(w) = \exp \left[ -\frac{\log \{p(w)\}}{\sum_{d=1}^D \sum_{j=1}^V n^{jd}} \right]$$

## [ Extracted Topic ]

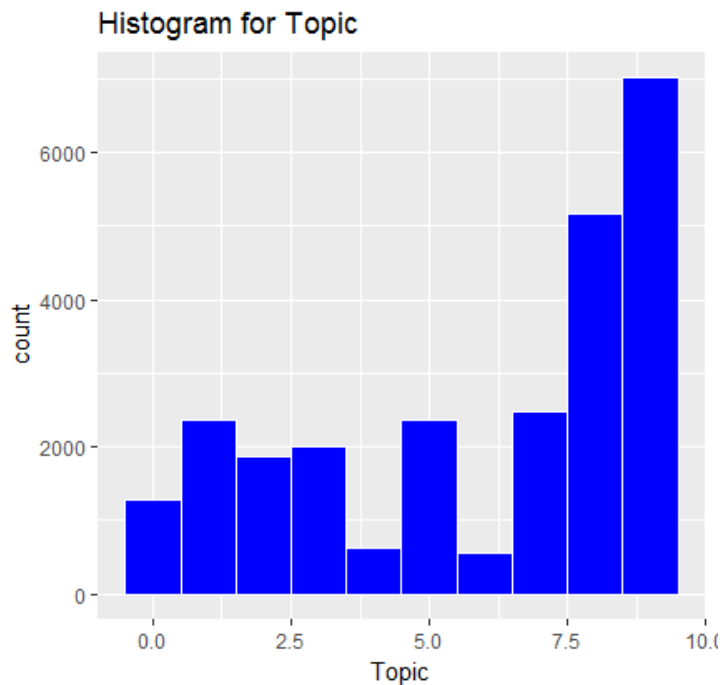
```
In [241]: pprint(ldamodel.print_topics())
```

```
[(0,
  '0.044*"'엄마"' + 0.032*"'돈"' + 0.028*"'아빠"' + 0.025*"'부모님"' + 0.017*"'아버지"' + 0.017*"'말"' +
  '+ 0.012*"'어머니"' + 0.011*"'가족"' + 0.010*"'동생"' + 0.010*"'생각"''),
 (1,
  '0.030*"'선배"' + 0.022*"'학교"' + 0.021*"'동아리"' + 0.021*"'수업"' + 0.020*"'학생"' +
  '0.016*"'교수님"' + 0.015*"'새내기"' + 0.012*"'학번"' + 0.011*"'중앙"' + 0.010*"'경희"''),
 (2,
  '0.029*"'학교"' + 0.017*"'자리"' + 0.013*"'기숙사"' + 0.012*"'공부"' + 0.011*"'도서관"' +
  '0.010*"'버스"' + 0.009*"'수업"' + 0.009*"'곳"' + 0.008*"'소리"' + 0.008*"'담배"''),
 (3,
  '0.277*"' + 0.035*"'이야기"' + 0.022*"'해"' + 0.016*"'되"' + 0.015*"'대나무"' + 0.014*"'시"' +
  '0.013*"'일상"' + 0.013*"'해주"' + 0.011*"'하기"' + 0.011*"'넌"''),
 (4,
  '0.025*"'노래"' + 0.015*"'언니"' + 0.014*"'게임"' + 0.012*"'술"' + 0.011*"'영화"' + 0.010*"'머리"' +
  '+ 0.010*"'음악"' + 0.010*"'형"' + 0.009*"'이름"' + 0.008*"'중대"''),
 (5,
  '0.041*"'공부"' + 0.025*"'생각"' + 0.023*"'학교"' + 0.019*"'대학"' + 0.018*"'시험"' + 0.014*"'누나"' +
  '+ 0.012*"'학점"' + 0.011*"'고등학교"' + 0.011*"'성적"' + 0.010*"'생활"''),
 (6,
  '0.110*"'오빠"' + 0.053*"'남자친구"' + 0.047*"'그녀"' + 0.036*"'여자친구"' + 0.016*"'연락"' +
  '0.016*"'전화"' + 0.015*"'생일"' + 0.013*"'죽순"' + 0.010*"'선물"' + 0.009*"'크리스마스"''),
 (7,
  '0.026*"'생각"' + 0.015*"'글"' + 0.015*"'문제"' + 0.013*"'사람"' + 0.012*"'말"' + 0.008*"'의견"' +
  '0.006*"'사회"' + 0.006*"'이유"' + 0.006*"'여성"' + 0.005*"'학생들"''),
 (8,
  '0.098*"'친구"' + 0.037*"'남자"' + 0.037*"'여자"' + 0.034*"'말"' + 0.030*"'사람"' + 0.029*"'생각"' +
  '+ 0.018*"'연락"' + 0.016*"'연애"' + 0.011*"'애기"' + 0.009*"'고민"''),
 (9,
  '0.063*"'사람"' + 0.046*"'생각"' + 0.039*"'사랑"' + 0.039*"'말"' + 0.030*"'마음"' + 0.017*"'행복"' +
  '+ 0.011*"'모습"' + 0.010*"'감정"' + 0.009*"'미안"' + 0.009*"'속"'')]
```

---

## [ 시각화 – Word Cloud & Cluster ]

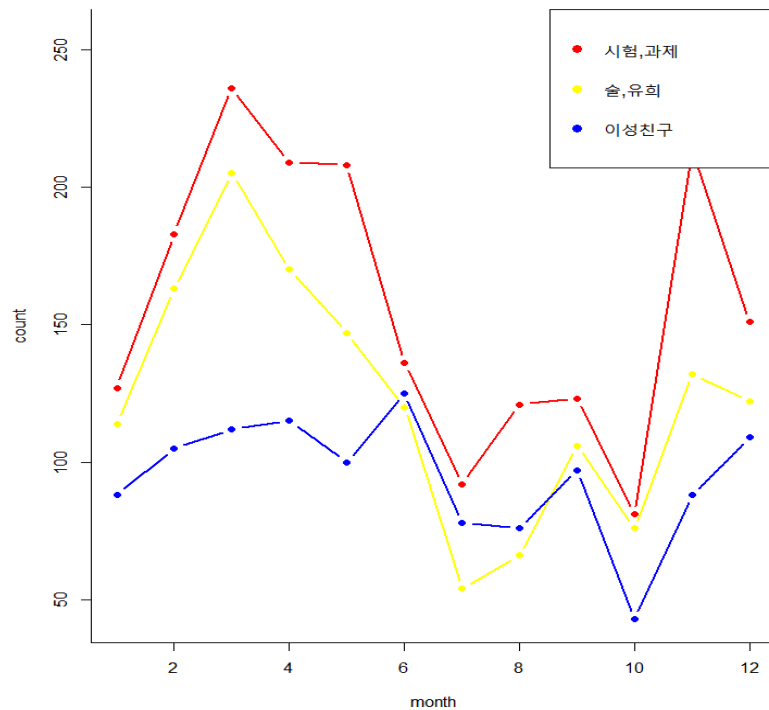
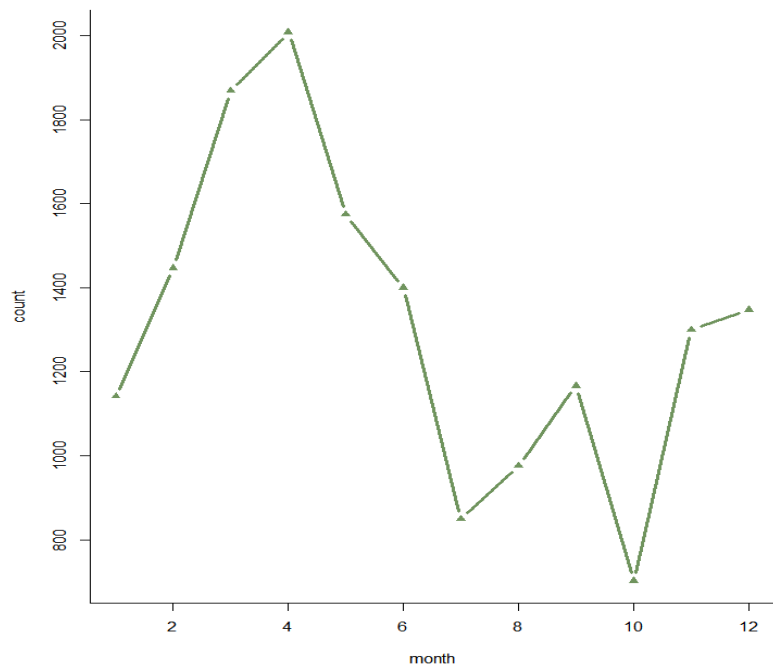
## [ Topic 상세 ]



### Topic 이름

- 0 : 가족
- 1 : 대학생 생활 (파생 생활)
- 2 : 시설
- 3 : 일상 (사는 얘기)
- 4 : 취미 / Entertainment
- 5 : 학업, 공부
- 6 : 연애 (썸)
- 7 : 사회 문제
- 8 : 친구 관계, 연애
- 9 : 사랑 (애인)

## [ 부록 - 월별 이용현황 ]



## [ 한계 - 아쉬운 점 $\pi^{\pi}$ ]

한계

- ! 끝없는 (전)전처리 지옥
- ! 대나무숲 빌런들의 글 전처리하기
- ! 말도 안되는 전처리 과정



KU BIG DATA

# TEXT MINING

“감사합니다”