

데이터분석캡스톤디자인

---

# 6주차 수행보고

**Khupid** 조

산업경영공학과 김동혁  
관광학과 류연주  
산업경영공학과 유정수

## 형태소 분석



## (2) 중요 문장 계산 후 요약

|    |                      |   |
|----|----------------------|---|
| 38 | 0.05179020833613878  | 그림전시장보다 훨씬 더 생생한 감이 있는 곳이 디자인 전시장 같습니다 .  |
| 13 | 0.05128247629951667  | 동시대 활동한 작가들의 작품도 일부 전시되어 있습니다 .   |
| 19 | 0.04251123001765623  | 전시장 가에 있는 철골에는 그가 디자인한 작품들에 대한 설명과 실물을 , 중앙에 있는 철골에는 붉은색 상자 안에 호기심이 많아 남의 집 서랍을 열어다 보는 듯한 느낌의 작품들이 전시되어 있어요 . |
| 3  | 0.039387673568825436 | 전시장 구성은 SPACE 1-5 까지 다섯 섹션으로 구분되어 있는데 , 구역마다 그의 작품의 특색이 명확하게 표현되어 있어 전시되어 있는 점이 좋았습니다 .                       |
| 0  | 0.03821818353286199  | 예술의 전당 한가람미술관에서 열리고 있는 이탈리아 디자인의 거장 카스틸리오니 전을 드디어 다녀왔습니다 .  |
| 22 | 0.03794562545912815  | 그도 뒤상과 같은 레디메이드 작품을 여럿 선보이고 있습니다 .  |
| 31 | 0.03742567135649523  | 34 점의 작품을 사방의 거울을 통하여 보니 정말 나무가 뻗어나 숲속에 있는 느낌이 들죠 ?   |
| 27 | 0.03421513498038972  | 포스터의 숲 이번 전시회가 특별한 두 가지 의미를 가지고 있는데 , 하나는 아시아 최초의 카스틸리오니 단독전이라는 점이고 .   |
| 23 | 0.033752482551147586 | 아이콘 이 전시실에서는 카스틸리오니 형제에게 국제적 명성을 안겨준 제품들이 전시되어 있습니다 .   |
| 34 | 0.032346722709293896 | 세상은 그의 작품을 놓고 ' 평범함의 승리 ' 라고 한다고 합니다 .  |
| 1  | 0.03173824892473695  | 이번 달 26 일까지이니 관심 있는 분 서두르세요 .   |
| 20 | 0.03084813364330009  | 눈 표시된 상자 안의 세계는 어쩌면 이번 전시장에서 상상으로 볼 수 있는 그의 건축물 디자인을 감상할 수 있는 코너였을 겁니다 .                                      |
| 18 | 0.030244442918920783 | 그래도 전시니 약간의 디자인 요소를 가미했는데 , 그것이 원치 않아맞혀 보라고 도슨트는 얘기합니다 .  |
| 다  | 0.02907994515746934  | 도슨트 채성도 뽐뽐한 마시구요  |

## (1) Okt (Open Korean text) 형태소 분석

```
In [6]: def get_tags_to_json(res):
        review_str = ''

        for review in review_df[review_df['res_name']==res]['review']:
            review_str+=str(review)

        okt = Okt()
        tags_dump = okt.pos(review_str, stem=True)
        tags = []

        for word in tags_dump:
            if word[1] == 'Adjective':
                tags.append(word[0])

        tags = '/'.join(tags)
        # text 파일로 저장
        with open('../data/review_tags/'+res+'.txt','w') as make_file:
            make_file.write(tags)
```

okt = Okt()  
tags\_dump = okt.pos(review\_str, stem=True)

## 형용사 추출

### <맞집>

'건장하다',  
'두드러지다',  
'무관심하다',  
'서늘하다',  
'처참하다',  
'흥사하다',  
'소란하다',  
'필요없다',  
'아담하다',  
'아박하다',  
'쉽다',  
'독하다',  
'음뻑하다',  
'적당하다',  
'문잡하다',  
'얕다',  
'넉넉하다',  
'토실하다',

### <전시>

'다르다',  
'이렇다',  
'궁금하다',  
'좋다',  
'뵤다',  
'많다',  
'새롭다',  
'소소하다',  
'어떨다',  
'아하다',  
'선하다',  
'특별하다',  
'뻥뻥하다',  
'평범하다',  
'자세하다',  
'평범하다',  
'생생하다',  
'평범하다',

## 중요 문장에서 형용사 추출

(( '좋았습니다', 'Adjective' ),) 0.193432  
(( '두꺼운', 'Adjective' ),) 0.193432  
(( '좋았습니다', 'Adjective' ),) 0.193432  
(( '같으면', 'Adjective' ),) 0.193432  
(( '어려웠을', 'Adjective' ),) 0.193432

# 불용어 처리 & 영어 번역 시도

## (1) 한글 불용어 사전 구축

Korean Stopwords

Home > Resources > Stopwords > Korean

**Korean Stopwords**

|        |        |           |
|--------|--------|-----------|
| 아      | 어찌됐든   | 하기보다는     |
| 휴      | 그위에    | 차라리       |
| 아이구    | 게다가    | 하는 편이 낫다  |
| 아이쿠    | 점에서 보아 | 흐흐        |
| 아이고    | 비추어 보아 | 놀라다       |
| 어      | 고려하면   | 상대적으로 말하자 |
| 나      | 하게될것이다 | 면         |
| 우리     | 일것이다   | 마치        |
| 저희     | 비교적    | 아니라면      |
| 따라     | 좀      | 쉴         |
| 의해     | 보다더    | 그렇지 않으면   |
| 을      | 비하면    | 그렇지 않다면   |
| 를      | 시키다    | 안 그러면     |
| 에      | 하게하다   | 아니었다면     |
| 의      | 할만하다   | 하든지       |
| 가      | 의해서    | 아니면       |
| 으로     | 연이서    | 이러면       |
| 로      | 이어서    | 좋아        |
| 에게     | 잇따라    | 알았어       |
| 뿐이다    | 뒤따라    | 하는것도      |
| 외가하여   | 뒤이어    | 그만이다      |
| 근거하여   | 결국     | 어쩔수 없다    |
| 입각하여   | 의지하여   | 하나        |
| 기준으로   | 기대어    | 일         |
| 예하면    | 통하여    | 일반적으로     |
| 예를 들면  | 자마자    | 일단        |
| 예를 들자면 | 더욱더    | 한편으로는     |
| 저      | 불구하고   | 오자마자      |

jupyter stop\_words2.txt 17시간 전

File Edit View Language

```
1 이 --- VCP -- 0.018279601
2 있 --- VA -- 0.011699048
3 하 --- VV -- 0.009773658
4 것 --- NNB -- 0.00973315
5 들 --- XSN -- 0.00689824
6 그 --- MM -- 0.005327252
7 되 --- VV -- 0.00361335
8 수 --- NNB -- 0.003473622
9 이 --- NP -- 0.003361203
10 보 --- VX -- 0.003310379
11 알 --- VX -- 0.0029757
```

```
In [190]: stop_words
Out[190]: {'그만이다',
            '하게하다',
            '연이서',
            '질마는지',
            '간지에서',
            '할 줄 안다',
            '이것',
            '의해',
            '시간',
            '정도',
            '라 해도',
            '그렇지만',
            '뒤이상',
            '시각',
            '어느 년도',
            '고로',
            '반대로 말하자면',
            '위대하면',
            '아하',
            '.....'}

In [191]: len(stop_words)
Out[191]: 664

In [192]: df = pd.DataFrame({'stop_word': stop_words})
           df.to_csv('all_stop_words.csv', index = False, header=True)
```

## (2) 한글 -> 영어 (papago 번역)

| review_en   | adv_en   |
|---|--|
| It's the best. It's comfortable. He's gone. Th... | [best, comfortable, only, favorite]                |
| Heesuya Hyosun king pork cutlet Ssangmun No.1     | []   |
| It was a little duhol. It's a little weak. Don... | [little, little, weak, first, wrong, female, o...] |
| A place where you can drink and enjoy bowling ... | [delicious, good, cheap, better, good, many, c...] |
| If you find a way to walk, you will be guided ... | [prettier, better, best, beautiful]                |
| ...   | ...  |
| There are a lot of fun and diverse board games... | [diverse, great]                                   |
| The cattle garden in the city center...Hidden ... | [best]   |
| It was clean and nice on the weekend! Chicken...  | [clean, nice, delicious, many, best, clean, de...] |
| NaN   | []   |
| Older people, dirty because they have a lot of... | [Older]  |

## (3) Tokenizing -> 형용사 추출

```
('clean', 199)
('comic', 98)
('bad', 90)
('quiet', 78)
('new', 71)
('comfortable', 70)
('cheap', 70)
('expensive', 67)
('hard', 62)
('neat', 61)
```

## Fasttext를 이용한 단어 벡터화

**pre-trained 한국어 모델**

```
In [9]: ko_model = models.fasttext.load_facebook_model('cc.ko.300.bin.gz')
```

```
In [11]: ko_model.wv.word_vec('문치있는')
```

```
Out[11]: array([ 3.23515316e-03, -1.44122355e-03, -8.17747638e-02,  4.61231880e-02,
        6.49930760e-02, -1.86859816e-03,  3.08580901e-02,  3.63530256e-02,
        4.65812832e-02,  3.09575372e-03, -4.47942875e-02, -1.39526548e-02,
       -9.68437642e-04,  5.35678901e-02,  1.20583080e-01,  5.72555773e-02,
        3.45472582e-02, -2.87518445e-02,  1.74929295e-02, -2.02273875e-02,
       -4.26698886e-02,  2.56524589e-02,  1.53578512e-04,  1.47120044e-01,
       -9.23347250e-02,  2.25054383e-01,  1.42958090e-02, -3.28313969e-02,
        7.17691844e-03, -3.52043323e-02, -8.20247531e-02,  9.34616756e-03,
       -1.27324024e-02, -3.66074927e-02, -1.87187623e-02, -8.31611380e-02,
        3.57123390e-02, -1.97619572e-02,  3.68778743e-02,  9.15034022e-03,
       -2.39420845e-03,  1.74316410e-02, -3.91573645e-02, -1.21175172e-02,
       -1.50920032e-02,  1.48614153e-01, -2.31579691e-02, -6.09047071e-04,
        2.43108217e-02, -5.14785685e-02, -2.24889815e-02,  7.58593110e-03,
       -9.17353705e-02,  3.79088591e-03,  8.14315006e-02, -5.26090898e-02,
        4.30432940e-03, -1.62116084e-02, -1.06273042e-02, -1.22643486e-02,
        4.80931215e-02, -2.30879411e-02, -1.24667808e-02, -7.11133704e-02,
        9.66753662e-02,  4.45196740e-02,  6.21299148e-02,  5.14498651e-02,
       -5.90744428e-02, -6.24508969e-02,  3.96463200e-02,  5.72216101e-02,
       -8.12430158e-02, -5.77316090e-02, -1.04852393e-02,  5.44773452e-02,
```

```
In [46]: ko_model.wv.word_vec('문치있는').shape
```

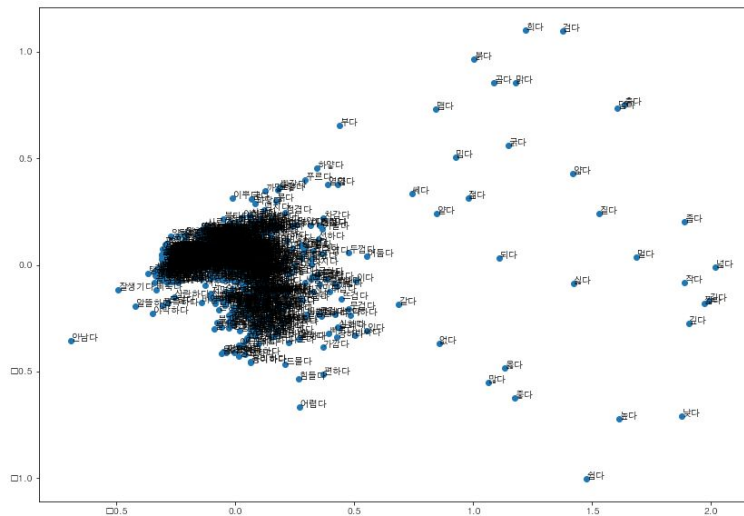
```
Out[46]: (300,)
```

PCA로 2차원으로 줄인 후 시각화한 모습

```
tags_vec.head()
```

| tag     | 0         | 1         | 2         | 3         | 4         | 5         | 6         | 7         | 8 ...         | 290      | 291       | 292       | 293       |       |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|---------------|----------|-----------|-----------|-----------|-------|
| 1 건강하다  | -0.033358 | -0.107392 | 0.079212  | -0.001186 | -0.048449 | -0.030069 | -0.007499 | 0.001004  | 0.133595 ...  | 0.013249 | 0.028847  | 0.045904  | -0.028165 | 0.13  |
| 2 두드러지다 | 0.018185  | 0.010256  | -0.047716 | -0.012996 | 0.004557  | -0.013141 | -0.018327 | -0.006541 | 0.017292 ...  | 0.010328 | -0.005689 | 0.014510  | 0.014435  | -0.00 |
| 3 무관하다  | 0.001222  | 0.031817  | 0.083566  | 0.017415  | -0.005664 | -0.016022 | -0.024136 | -0.002087 | -0.029156 ... | 0.014161 | 0.028465  | 0.032042  | -0.075265 | -0.01 |
| 4 서둘다   | -0.035465 | 0.106854  | 0.156820  | -0.054557 | 0.041650  | 0.000100  | -0.079635 | -0.158097 | -0.081634 ... | 0.072447 | -0.029971 | 0.035325  | -0.101966 | -0.05 |
| 5 저참하다  | -0.076719 | 0.009150  | 0.046816  | -0.033846 | 0.011217  | -0.013326 | 0.062515  | -0.017747 | -0.009819 ... | 0.055262 | -0.029961 | -0.042899 | 0.036459  | 0.00  |

5 rows x 301 columns



# 어떤 모델을 어떻게

쓸지?

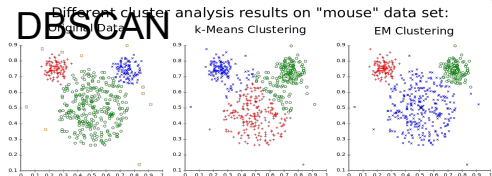
1) 한국어 리뷰 -> 영어로 바꿔서도  
지 해

| review   | review_en   |
|--|---|
| 최고예요. 편안하게 있다 갔어요. 분위기<br>기 작살!! 분위기조야요 연남동에서<br>유일하게... | It's the best. It's<br>comfortable. He's<br>gone. Th... |
| 희수야호선왕따스 쌍문1줄번게  | Heesuya Hyosun<br>king pork cutlet<br>Ssangmun No.1     |

2) 단어 벡터화 - FastText

```
Out[11]: array([ 3.23515316e-03, -1.44122355e-03, -8.17747638e-02,  4.61231880e-02,  
 6.49930760e-02, -1.86859816e-03,  3.08580901e-02,  3.63530256e-02,  
 4.65812832e-02,  3.09575372e-03, -4.47942875e-02, -1.39526548e-02,  
 -9.68437642e-04,  5.35678901e-02,  1.20583080e-01,  5.72555773e-02,  
 3.45472582e-02, -2.87518445e-02,  1.74929295e-02, -2.02273875e-02])
```

3) 클러스터링 - Kmeans, GMM,  
DBSCAN



# 다음주 계획

1) 전처리 마무리 및 필드 통일 : 동혁

2) 클러스터링 결과 최적화 : 연주, 정수

3) 클러스터링 기반으로 매트릭스 생성

DataSet5 <형용사 태그 데이터셋>

| ID | 떠들썩한 | 조용한 | 운치있는 | 가족끼리 | ... |
|----|------|-----|------|------|-----|
|    | 1    | 0   | 0    | 1    |     |
|    | 0    | 1   | 1    | 0    |     |
|    | 0    | 1   | 1    | 1    |     |