

캡스톤 디자인(2) 최종 발 표

Web Crawling, NLP를 활용한
Word Cloud 기반 기사 통합 요약 시스템

목차

1. 프로젝트 소개
 - 개발 동기, 개요, 배경
2. 설계 및 구현
 - Software architecture, 전처리 결과, 구현 내용
3. 실행 결과
 - 시연 및 설명
4. 향후 확장 가능성
5. 참고 문헌



프로젝트 소개

프로젝트 개발 동기



프로젝트 개발 동기



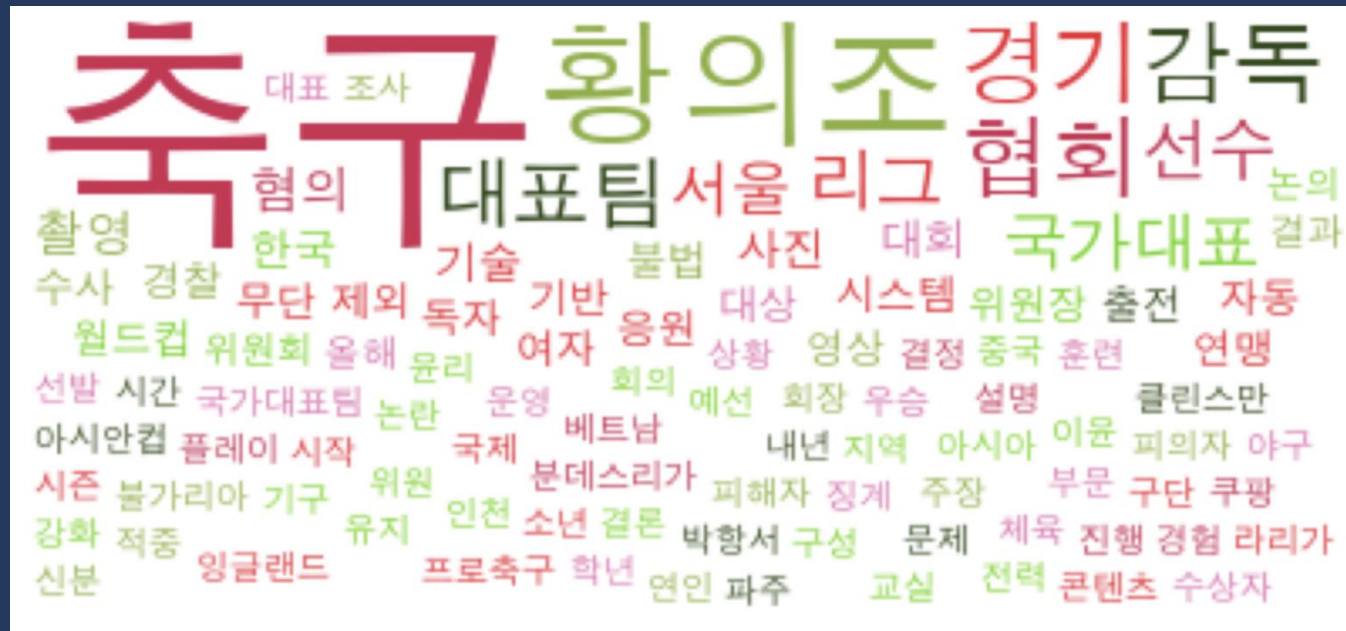
프로젝트 개요



프로젝트 개요



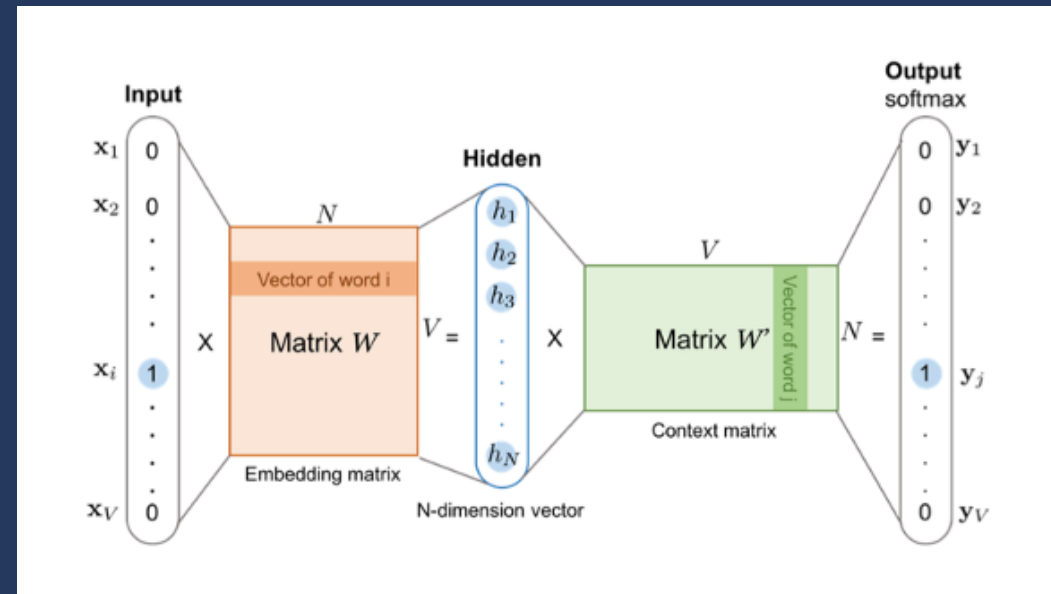
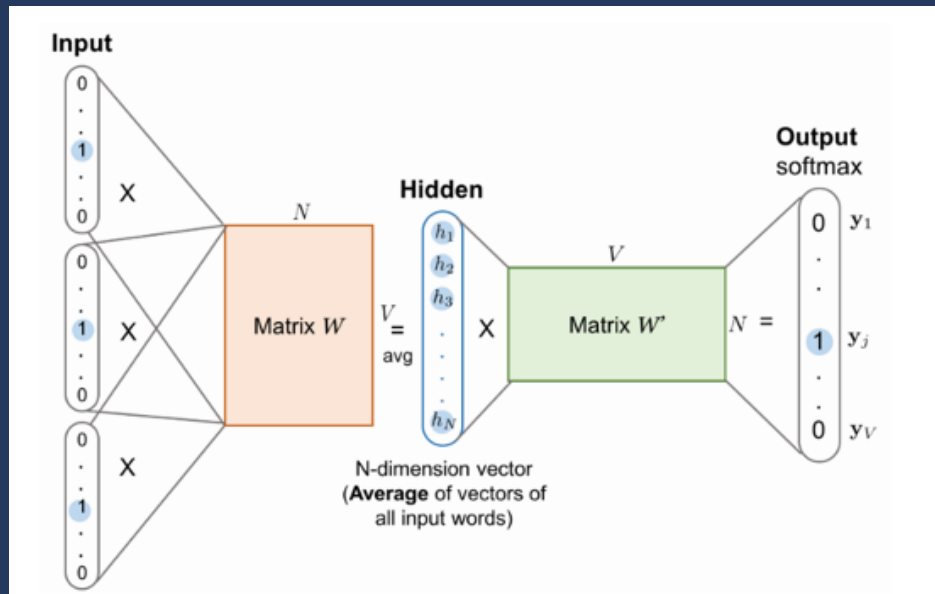
프로젝트 개요



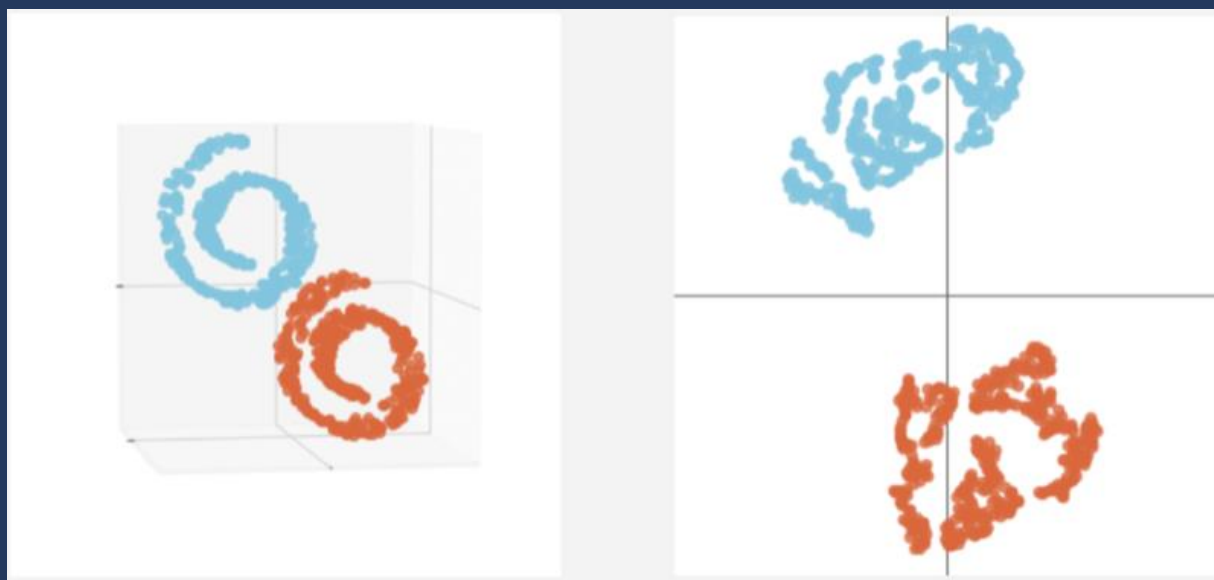
프로젝트 배경 : KoNLPy



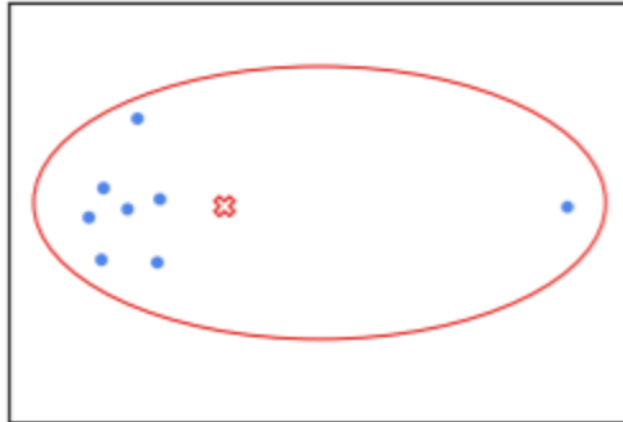
프로젝트 배경 : Word2Vec



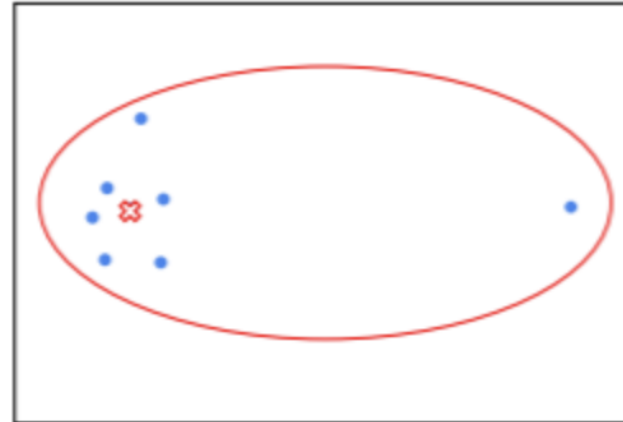
프로젝트 배경 : t-SNE



프로젝트 배경 : K-medoids

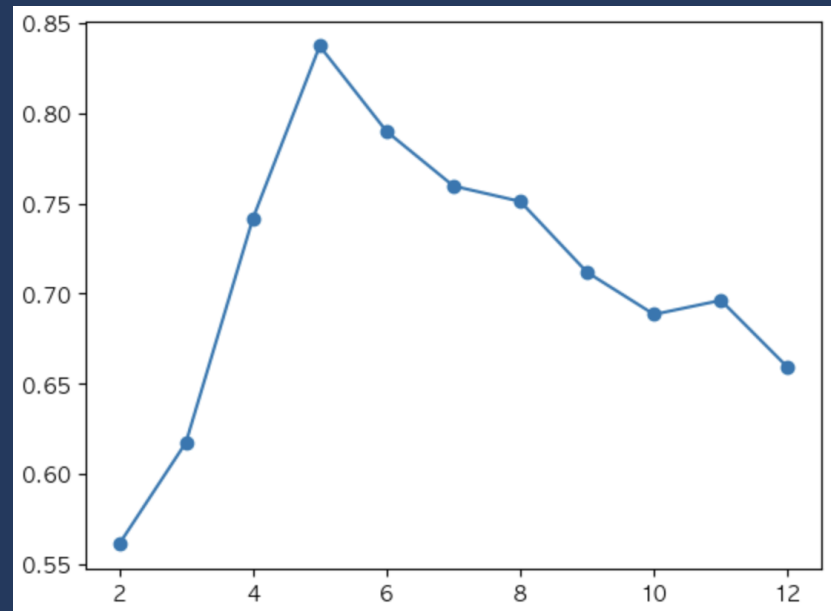


(a) Mean

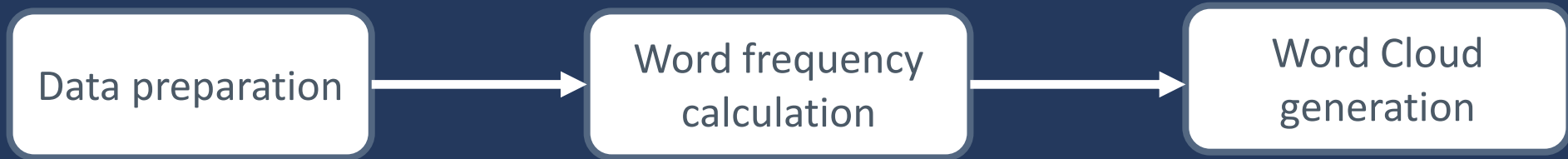


(b) Medoid

프로젝트 배경 : K-medoid – Silhouette method



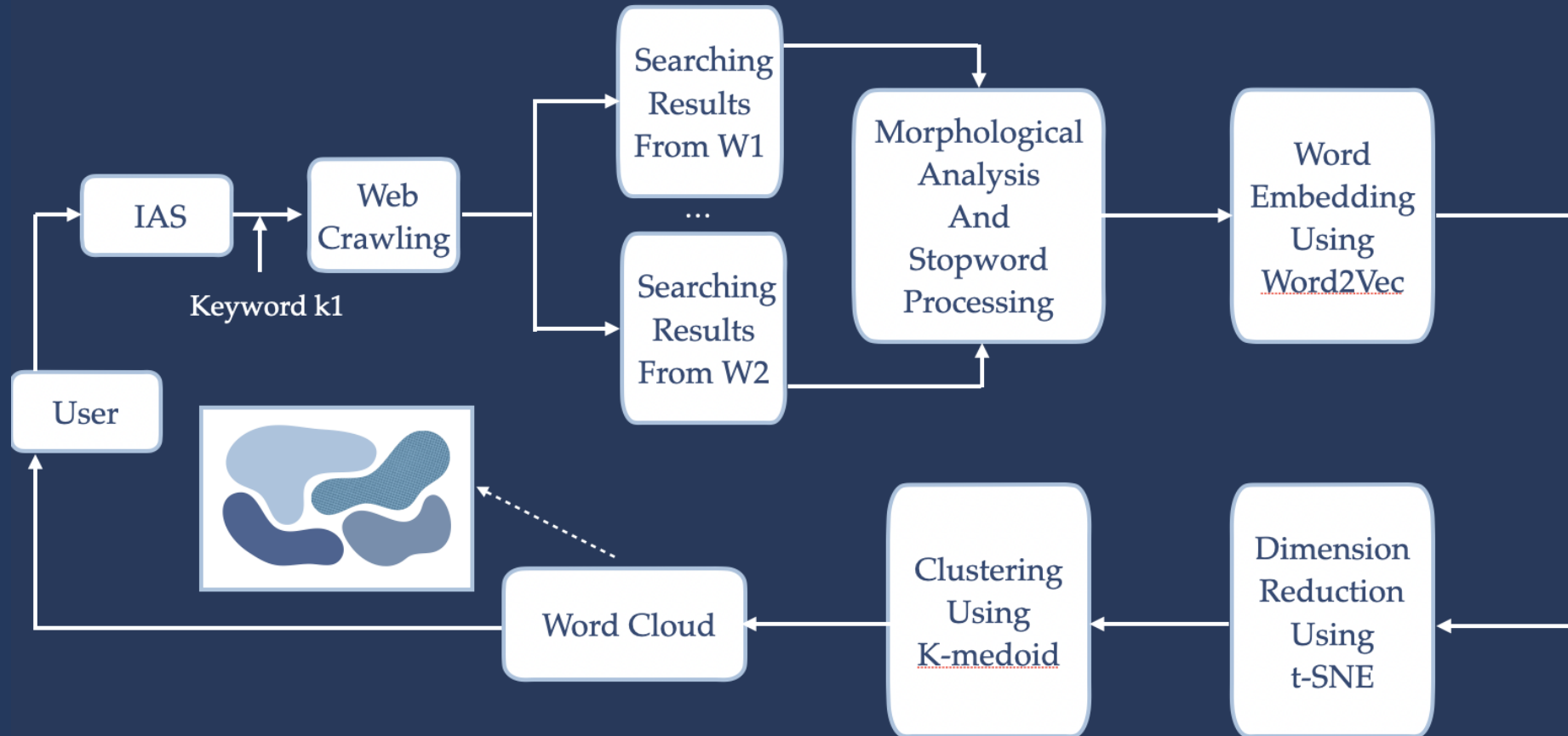
프로젝트 배경 : Word Cloud





설계 및 구현

Software Architecture



Dataset: Web Crawling

	date	link	content
0	2023.11.28.	https://n.news.naver.com/mnews/article/001/001...	[축구 국가대표 황의조[연합뉴스 자료사진](서울=연합뉴스) 안홍석 기자 = 축구 국...
1	2023.11.28.	https://n.news.naver.com/mnews/article/366/000...	[대한민국 축구 국가대표팀 황의조 /뉴스1 축구 국가대표 황의조(노리치...
2	2023.11.28.	https://n.news.naver.com/mnews/article/052/000...	[©연합뉴스불법 촬영 혐의로 경찰 수사를 받는 축구 국가대표 황의조(31·노리치시티...
3	2023.11.28.	https://n.news.naver.com/mnews/article/020/000...	[대한민국 축구 국가대표팀 황의조가 19일 오전 2026 FIFA 북중미 월드컵 아...
4	2023.11.28.	https://sports.news.naver.com/news.nhn?oid=421...	[(서울=뉴스1) 김성진 기자 = 이윤남 대한축구협회 윤리위원장이 28일 오후 서울...

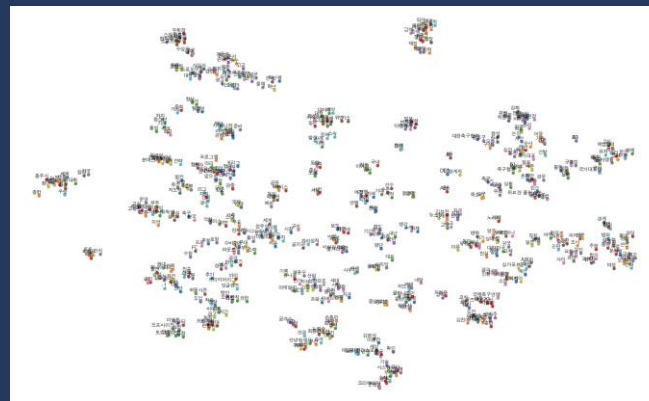
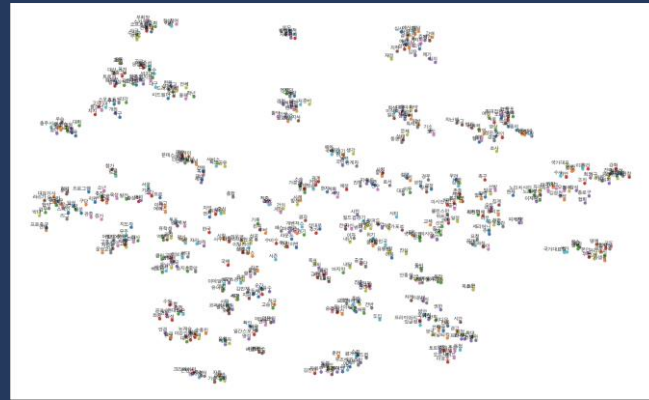
Preprocessing : 형태소 분석 및 불용어 처리



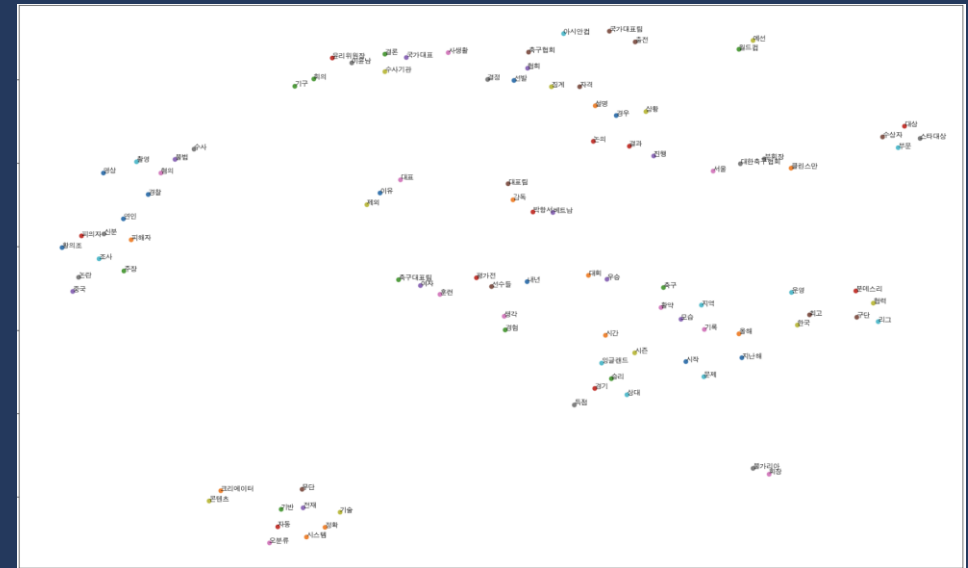
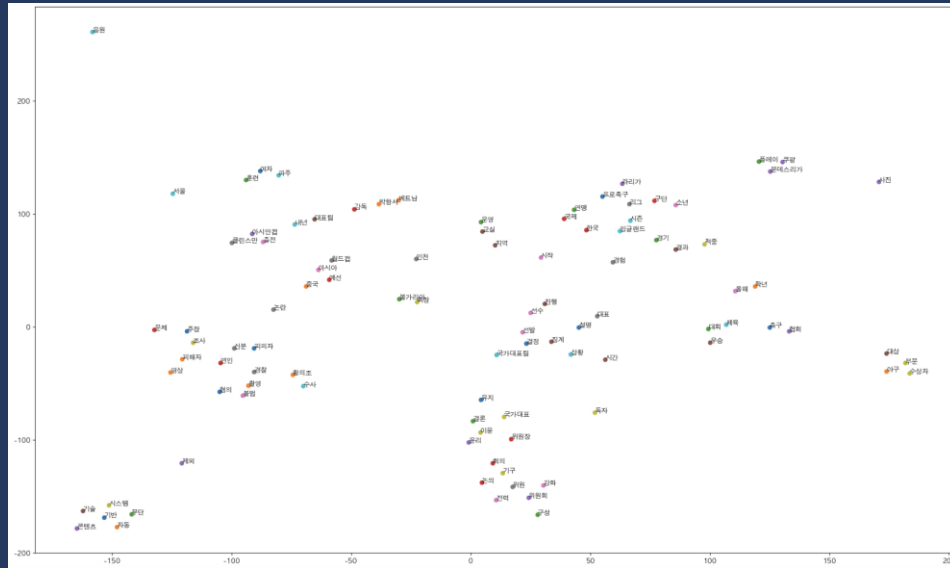
Preprocessing : Word2Vec

```
model_kkma = word2vec.Word2Vec(sentences=preprocess_kkma, vector_size=300, window=10, min_count=5, sg=1)
model_okt = word2vec.Word2Vec(sentences=preprocess_okt, vector_size=300, window=10, min_count=5, sg=1)
model_hannanum = word2vec.Word2Vec(sentences=preprocess_hannanum, vector_size=300, window=10, min_count=5, sg=1)
model_komoran = word2vec.Word2Vec(sentences=preprocess_komoran, vector_size=300, window=10, min_count=5, sg=1)
```

Preprocessing : t-SNE 기반 차원 축소



Preprocessing : t-SNE 기반 차원 축소 및 빈도 수 포함



Preprocessing : t-SNE 기반 차원 축소 및 빈도 수 포함

```
print(labels_okt[30])
```

여자

```
print(labels_okt[-1])
```

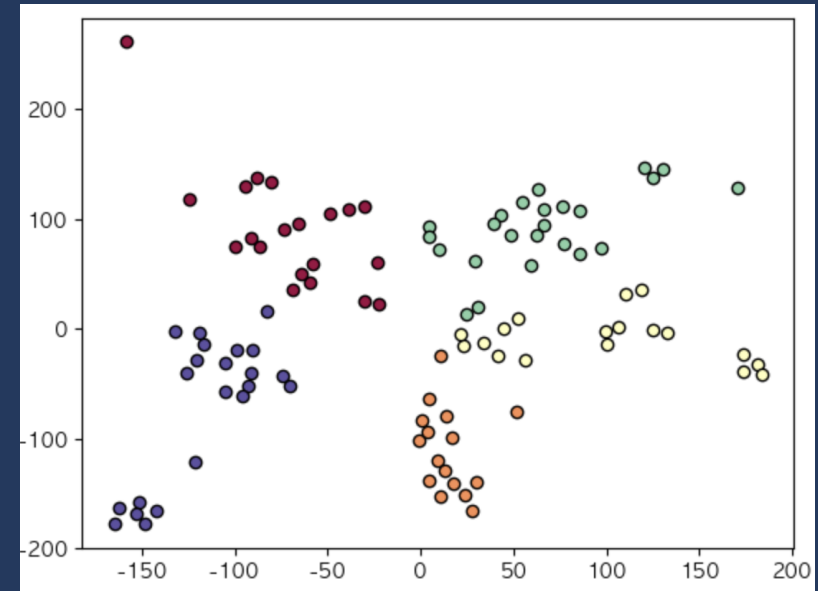
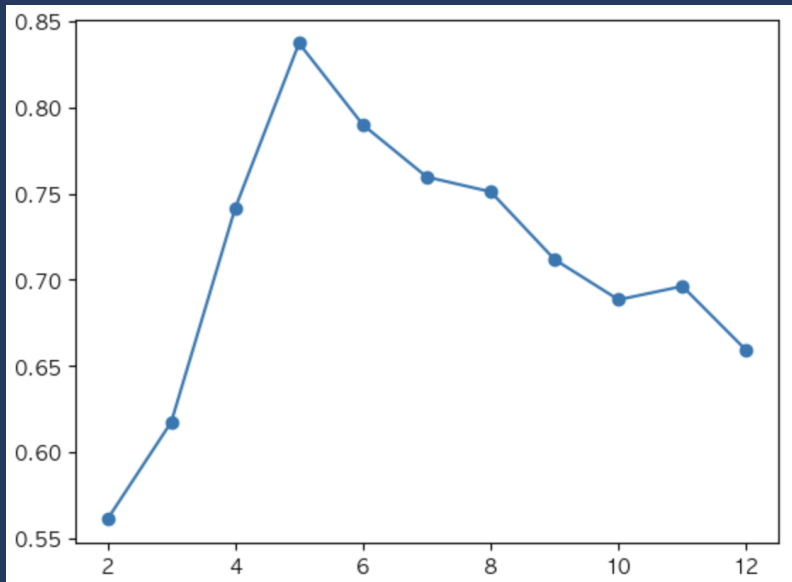
파주



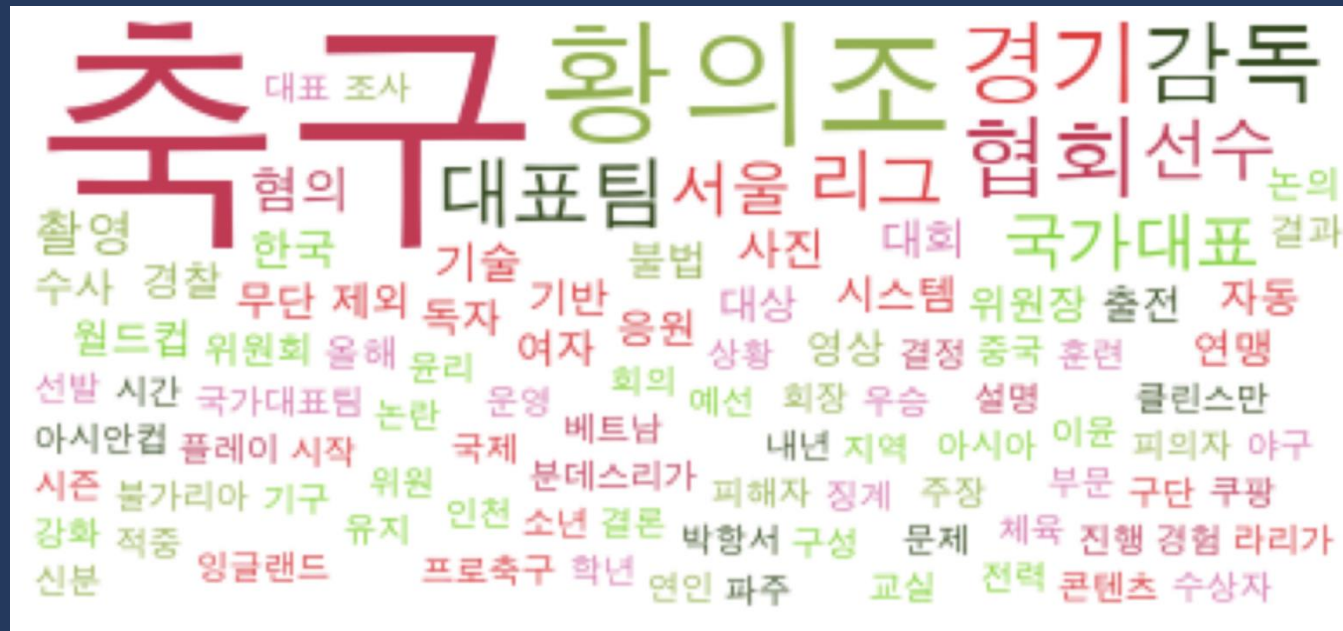
```
print(cos_sim(values_okt[30], values_okt[-1]))
```

0.99949986

Preprocessing : K-medoids



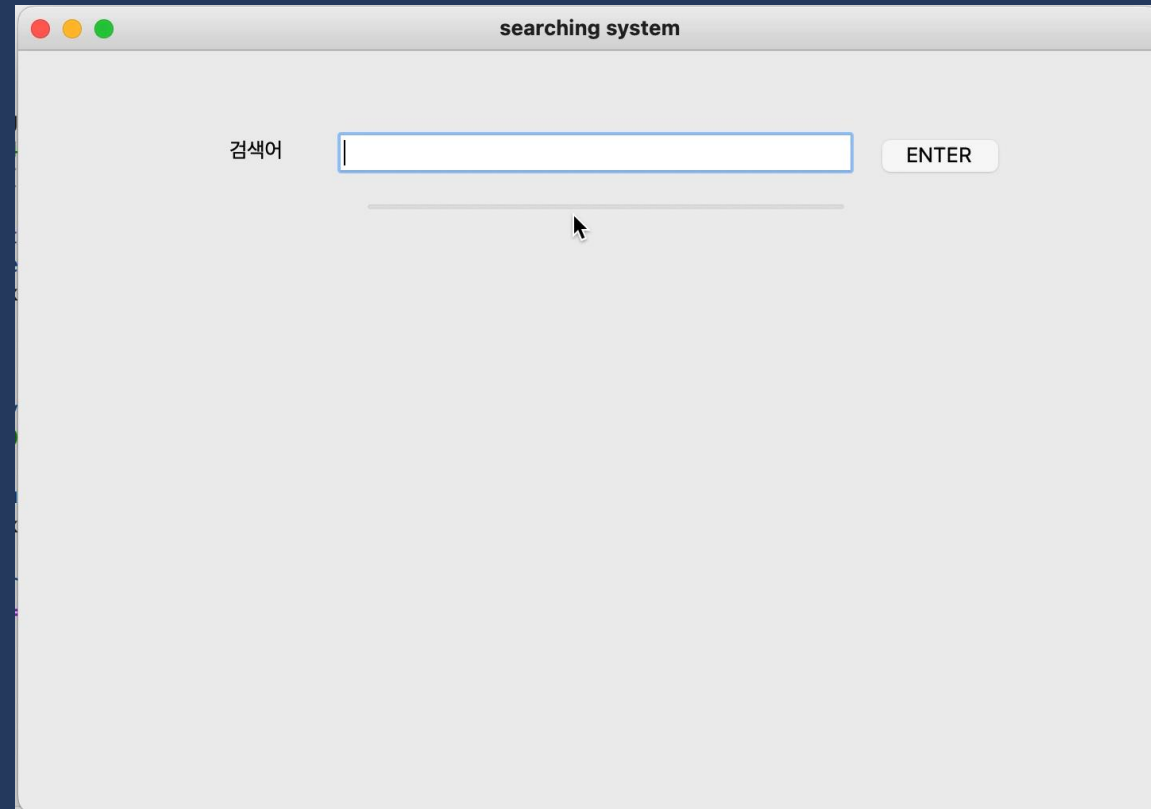
Word Cloud : 클러스터링 기반 색상 및 재배치





실행 결과

시연 결과





향후 확장 가능성

프로젝트 의미와 향후 확장 가능성

- 정보 습득의 효율 상승
- Word Cloud의 시각화 기능 향상
- Sematic Cloud의 대중화
- 이후 병렬화 도입하여 사용 가능성 향상 예정

참고 문헌

참고 문헌

- ‘Natural Language Processing with Transformers; Building Language Applications with Hugging Face’, Lewis Tunstall, Leandro von Werra & Thomas Wolf, O’Reilly
- ‘Word Cloud Explorer: Text Analytics Based on Word Clouds’, Heimerl, Florian. 2014 47th Hawaii International Conference on System Sciences ISBN
- ‘Semantic Word Cloud Generation Based on Word Embeddings’, Jin Xu, Yubo Tao, Hai Lin
- ‘Context Preserving Dynamic Word Cloud Visualization’, Weiwei Cui, Yingcai Wu, Shixia Liu, Furu Wei, Michelle X.Zhou