

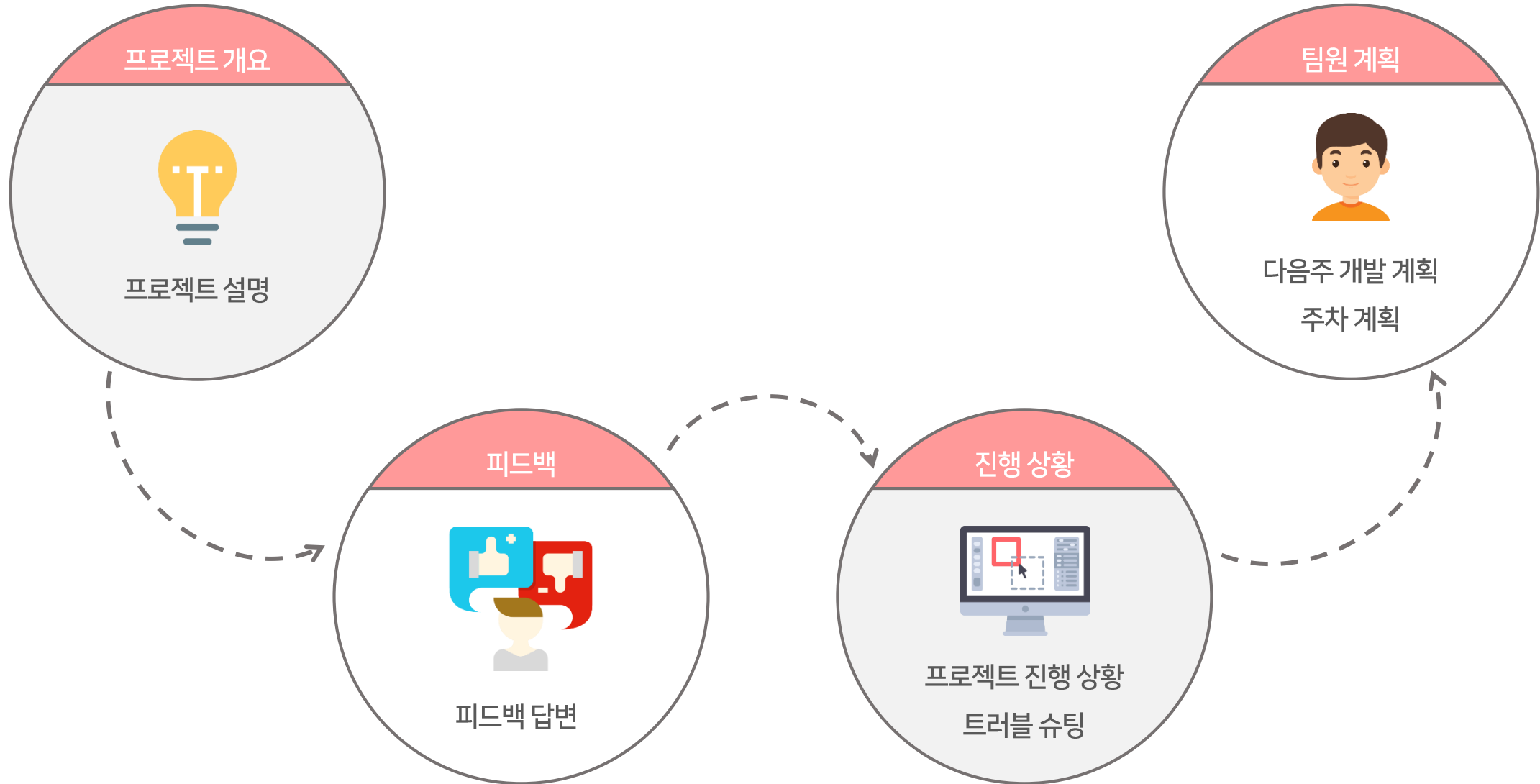


# 커뮤니티 상 거래를 위한 게시글 제목 변경 웹 확장 프로그램



팀명 : 제뿔  
박시현 + 박진영 + 이민희

# TABLE of CONTENTS

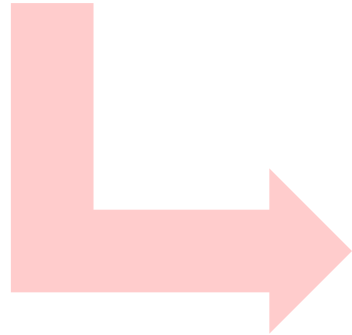


## 프로젝트 설명

♥ 무료배송 ♥ 워싱 면 라이더 자켓 판매 🧑🏻 N

니트, 필립 림 자켓, 양가죽 패딩코트(BCBG, Part or lady, 마이클 코어) 판매 🧑🏻 N

\*\* ♠️ ♠️ 셀린느 지고트 레니본 시스템 CK 존스뉴욕 아디다스 모르간 MNG ♠️ ♠️ 4  
판매 🧑🏻 [2] N



워싱 면 라이더 자켓 | 15000 | 사용감 유 | 2주

필립 림 자켓 | 30000 | 새상품 | 0회

존스뉴욕 | 43000 | 최상 | 착용만

## 지난 주 피드백

---

Bert 사용 고려?



중고나라 데이터 특성상 문맥의 흐름에 의존하지 않고도  
단어의 특징을 잘 반환할 것이라고 생각

## 프로젝트 진행 상황

박시현

데이터셋 추가 구축

서버에 크롤링 프로그램 등록

박진영

텍스트 대치 테스트 데이터로 구현

웹페이지 데이터 저장

이민희

NER 모델 구축

모델 파라미터 변경

## 모델 학습 위한 데이터셋 추가 확보

게시글 오류 / 사용자 탈퇴 등 데이터 수집 오류 코드 수정

[노트북, 남성상의]

2개 카테고리 별로 각각 10000개씩 추가 확보

## 게시물 다양성 부족

최대 목록 페이지 = 1000페이지

1000번째 페이지 게시물들의 작성일은 2019.11.01 (여성 상의 게시판 기준)

의류의 경우 계절별 의류 종류가 달라 문제



문제점을 인지하고 획득 가능한 데이터로만 학습

# 우분투 서버 연결

---

우분투 서버 구축

필요한 모듈 설치 및 환경설정



```
ubuntu@ip-172-31-34-225:~/capstone$ ls  
naverlogin.py  requirements.txt  sites  utils
```



## DOMSubtreeModified

```
var highestScore;

console.log('autoEvent listener!');
document.addEventListener('DOMSubtreeModified', function() {
  console.log('DOMSubtreeModified listener!');
  if(document.body.innerHTML.search('main-area') != -1) {
    console.log('target frame!');

    var table = document.getElementsByClassName('article-board');
    console.log(table.getElementsByTagName('a'));

    var titleList = table.getElementsByTagName('a');
```

DOMSubtreeModified가 필요한지 깊이 고려 X

## DOMSubtreeModified

```
2 ▶ Uncaught RangeError: Maximum call stack size exceeded.  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)  
   at HTMLDocument.<anonymous> (auto.js:26)
```

3

4

DOMSubtreeModified 이벤트 리스너 내부에서 innerText 사용  
→ 무한루프 발생

반복 횟수 제한 → 호출 스택 순서 실행으로 결과 중복  
→ DOMSubtreeModified 이벤트 리스너 제거 (문제 해결)

## 텍스트 대치 구현 상황

	제목	작성자	작성일	조회
654912086	제목a 판매 🧑	자입지	2019.10.19.	17
646926413	제목b 판매 🧑	positiveart	2019.10.01.	9
642716194	제목c 판매 🧑	한옥부부	2019.09.21.	35
626446408	제목b 판매 🧑 [1]	SONGEI0901	2019.08.11.	19
623635707	제목e 판매 🧑	궁디애기	2019.08.04.	55
620970984	제목f 판매 안전 🧑	jinhee0448	2019.07.28.	85
618993083	제목g 판매 🧑	갯요정	2019.07.23.	18
616581663	제목h 판매 🧑	뉴앙스탈사랑스	2019.07.14.	37
615476511	제목i 판매 🧑	loveick	2019.07.14.	37
614007547	제목j 판매 🧑	loveick	2019.07.14.	37
606060821	제목k 판매 안전 🧑	p37625	2019.06.20.	22

확장 앱에서 변경 버튼 클릭  
→ 다음과 같이 변경



## 쿠키 vs 웹스토리지

### 쿠키

작은 데이터 저장  
서버와 통신하여 데이터 가져옴

### 웹스토리지

비교적 큰 데이터 저장  
서버 통신 불필요

→ 웹스토리지 사용

## utf-8 문제

[ Error: ' cp949 ' codec can ' t decode byte 0xec in position 0: illegal multibyte sequence ]

```
# -*- coding: utf-8 -*-  
from pprint import pprint  
from konlpy.tag import Twitter #customized twitter  
from konlpy.tag import Kkma #문장 분리용  
import gensim  
  
twitter = Twitter()  
  
def read_data(filename):  
    with open(filename, 'r', encoding='UTF8') as f:  
        data = f.read().splitlines()  
    return data
```

## kkma 문장 분리 라이브러리

```
10: '시 검 형 x1, 검 핑 m 입니다'
11: '직거래 산 본 또는 금정 역 택배거래 시 +5000x1 검
12: '1x 나이키 빅 스우'
13: '시 후 리스 검 형 ,m 검 핑 팝니다'
14: '[ 판매 ]이 얼 즈 어고 패 커 블 스웨트 셔츠 02 오렌...
15: '이 얼 즈 어고 패 커 블 스웨트 셔츠 02 판매 하보 43
16: '[ 판매 ]이 얼 즈 어고 패 커 블 스웨트 셔츠 02 오렌...
17: '코오롱 스노우 볼 카 키, 블랙 색상 상급 팝니다'
18: '그레이 색 카 키 블랙 색상 상급 팝니다'
```

지나친 띄어쓰기 및 문장분리 (토큰화에 불편)

→ kkma 라이브러리 대신 다른 라이브러리 사용

# 불필요한 띄어쓰기 제거

```
02: '상태 최상입니다'
03: '디자인, 핏 다 이빠요신랑한테 안어울려서 내놓아요'
04: '젊은분취향입니다'
05: '인천직거래)코모도스퀘어 울100% 세미오버핏 싱글코트'
06: '1x나이키 빅스우시 후리스 검형,m 검핑팝니다'
07: '홍대 나이키 새제품 나이키 빅스우시 검형 x1,검핑m입니다'
08: '직거래 산본또는금정역택배거래시+5000x1검형22m핑검19일괄 구매시 390000010 5261 6722'
09: '1x나이키 빅스우시 후리스 검형,m 검핑팝니다'
10: '[판매]이얼즈어고 패커블 스웻셔츠 02 오렌지 M'
11: '이얼즈어고 패커블 스웻셔츠 02 판매착불 13만9천원5회 실착쿨거래하실분 연락주세요'
12: '[판매]이얼즈어고 패커블 스웻셔츠 02 오렌지 M'
13: '코오롱 스노우볼 카키, 블랙색상 상급 팝니다'
14: '그레이쉬 카키, 블랙 색상 105사이즈입니다.'
15: '코트위주로 입다보니 착용횟수가굉장히 적습니다.'
```

```
13
14 def sentenceSlice(phrase_list):
15     sentence_list = []
16     for phrase in phrase_list:
17         sentence_list.extend(kss.split_sentences(phrase))
18     return sentence_list
```

kss 라이브러리 사용

(문장 분리 속도개선 + 불필요한 띄어쓰기 제거)



## 불필요한 띄어쓰기 제거

13: '시 후 리스 검 형 ,m 검 핑 팝니다'

14: '[ 판매 ]이 얼 즈 어고 패 커 블 스웨 셔츠 02 오렌...

15: '이 얼 즈 어고 패 커 블 스웨 셔츠 02 판매 차분 43  
[ 판매 ]이 얼 즈 어고 패 커 블 스웨 셔츠

16: '[ 판매 ]이 얼 즈 어고 패 커 블 스웨 셔츠 02 오렌...

17: '코오롱 스노우 볼 카 키, 블랙 색상 상급 팝니다'

09: '1x나이키 빅스우시 후리스 검형,m 검핑팝니다'

10: '[판매]이얼즈어고 패커블 스웨셔츠 02 오렌지 M'

11: '이얼즈어고 패커블 스웨셔츠 02 판매차분 13만9천원5회 실착클거래하실분 연락주세요'

12: '[판매]이얼즈어고 패커블 스웨셔츠 02 오렌지 M'

13: '코오롱 스노우볼 카키, 블랙색상 상급 팝니다'

14: '그레이쉬 카키, 블랙 색상 105사이즈입니다.'



## 시각화 및 정확도

다차원의 word2vec을 시각화를 위해 2차원으로 축소  
축소할 때 관계를 유지하기 위해 t-SNE 사용

t-SNE : t-distributed Stochastic Neighbor Embedding

고차원 공간에서의 유클리디안 거리측정방법 활용해  
데이터 포인트들의 유사성을 표현하는 조건부 확률로 변환하는 방법

## 시각화 및 정확도

단점 : 조건부 확률의 기준이 정해져 있지 않아  
생성할 때마다 모양이 다르다

- 대략적인 정확도 유추 가능하지만
- 정확한 정확도 지표로는 사용 불가능
- 새로운 방법 모색 중

## 토큰화 시 필요로 하는 품사 변경

알파벳, 숫자

similarity 비교 시 단어 벡터 학습이 잘 안되는 문제점

Similarity 처리 대신에 다른 처리 모색 중

→ 명사, 형용사, 동사만 토큰화 후 사용

## NER

```
[('푸마 /Noun', 0.5078113675117493),  
('소우 /Noun', 0.469389945268631),  
('리복 /Noun', 0.4629251956939697),  
('뉴발란스 /Noun', 0.4623837172985077),  
('루나 /Noun', 0.45955997705459595),  
('줌 /Noun', 0.45431721210479736),  
('아디다스 /Noun', 0.4525972008705139),  
('카모 /Noun', 0.4494244456291199),  
('콜라보 /Noun', 0.43964675068855286),  
('후디 /Noun', 0.4356365203857422),  
('데상트 /Noun', 0.4347517788410187),  
('휠라 /Noun', 0.4312279224395752),  
('카파 /Noun', 0.4295527935028076),  
('리액트 /Noun', 0.4268908202648163),  
('레오 /Noun', 0.4262371063232422),  
('에리트 /Noun', 0.4244637101305634)]
```

이름 브랜드를 처리해주는  
사전을 만들기 위한 모델 생성  
→ 나이키, 아디다스와 유사도가  
높게 추출된 목록 출력

# Logistic Regression

패커블 /Noun 16  
스노우볼 /Noun 10  
평창 /Noun 13  
모노크롬 /Noun 10  
케이스위스 /Noun 33  
곰돌이 /Noun 24  
빨릴레리 /Noun 20  
빅바 /Noun 20  
녹 /Noun 20  
단군 /Noun 14  
쉐펠 /Noun 20  
뱅뱅 /Noun 39  
에르메스 /Noun 42  
일삼일삼 /Noun 16  
레알마드리드 /Noun 11  
볼빅 /Noun 46  
휠르자 /Noun 14

Logistic Regression을 이용,  
위에서 추출한 seed words로  
positive class 예측

비슷하다고 예측된 단어가 선택

BUT, 브랜드 같지 않은 단어들 출력

# 정리

---

seed words → 사전에 추가

Logistic Regression 추출 값 → 튜닝 예정



# 정리

---

여성 옷 / 남성 옷 모델을 나눠서 생성 예정  
(옷 종류, 구매하는 브랜드의 차이 발생)

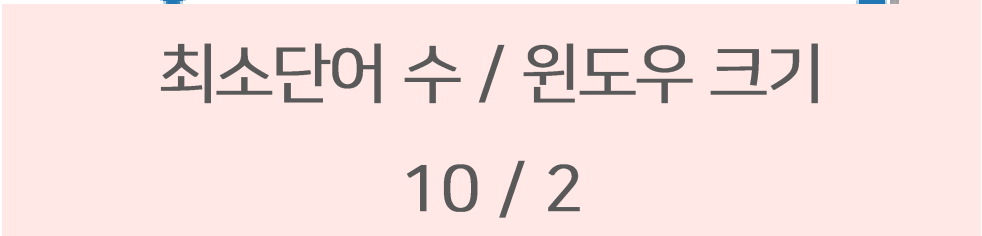
→ 모델을 각각 만들었을 때의 정확도 상승 예상

# 모델 파라미터 변경

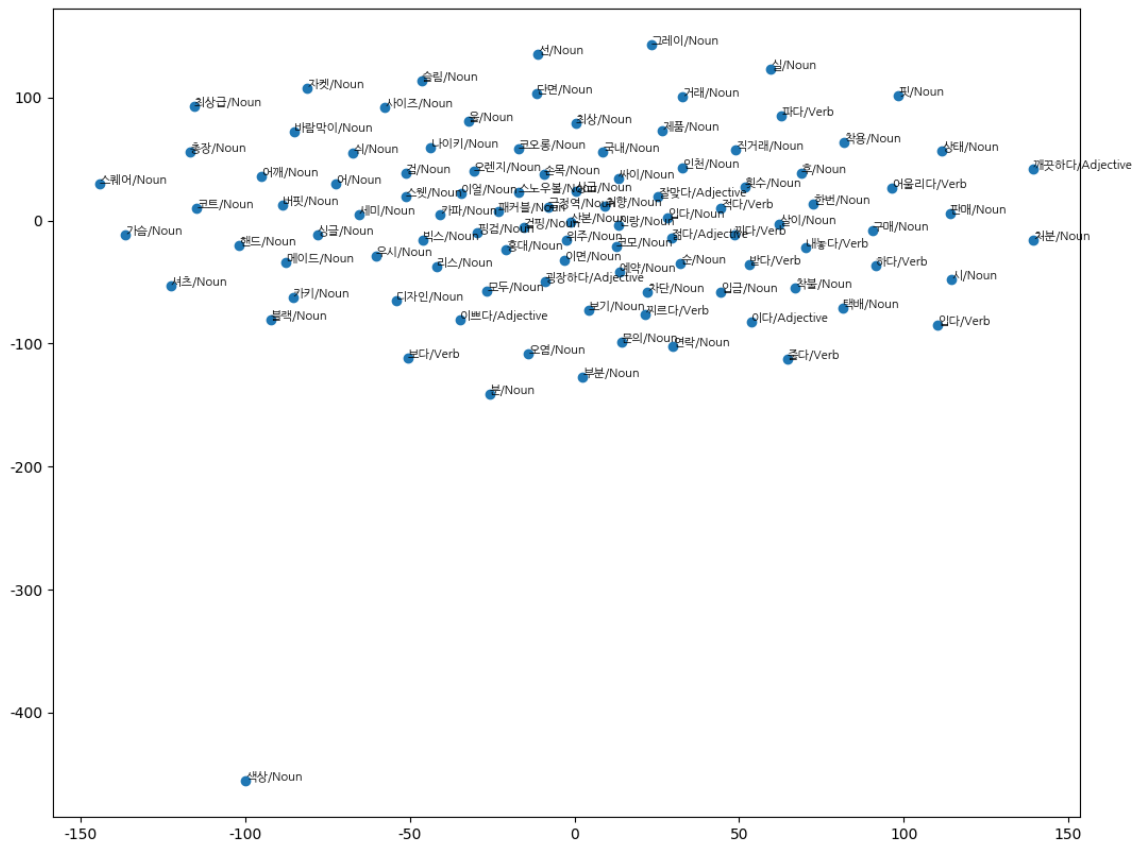
```
def makeM(train_data):  
    #하이퍼 파라미터 변경하는 작업을 함  
    num_features = 300    #문자 벡터 차원 수(100~300)크면 차원의 저주 최대 6천만 단어의 말뭉치에 300벡터  
    min_word_count = 5  
    #최소 문자 수(10~100) 최소 문자 빈도가 적으면 학습이 잘 안될 수 있음  
    #단어마다 문맥 학습 최소 빈도수가 다름(ex 트와이스는 쉬운데, 아프리카tv, 대륙아프리카가ㄷ이 여러 의미로 사용되는)  
    num_workers = 4    #병렬 처리 쓰레드 수(2or4)  
    context = 2    #문자열 창 크기 (2~10) 창 크기가 크면 의미적 결과 작으면 구문적 결과  
    downsampling = 1e-3    #문자 빈도 수 Downsample  
  
    model = gensim.models.Word2Vec(train_data, workers=num_workers,  
                                   size=num_features, min_count = min_word_count,  
                                   window = context, sample = downsampling)
```

여러 예제들 참고 후 값 설정





# 모델 비교



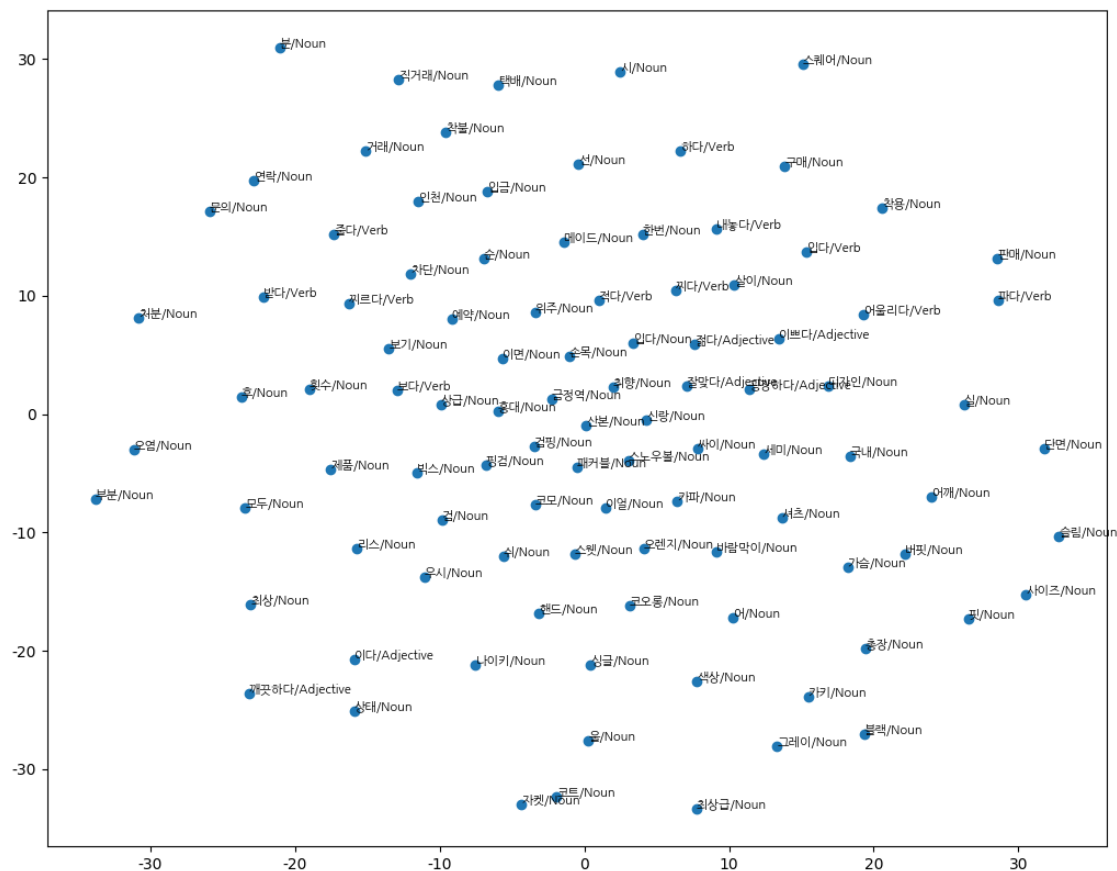
최소단어 수 / 윈도우 크기

10 / 6

하나의 예외값

토큰이 잘 안 나뉘져 있음

# 모델 비교



최소단어 수 / 윈도우 크기

10 / 4

잘 출력됨

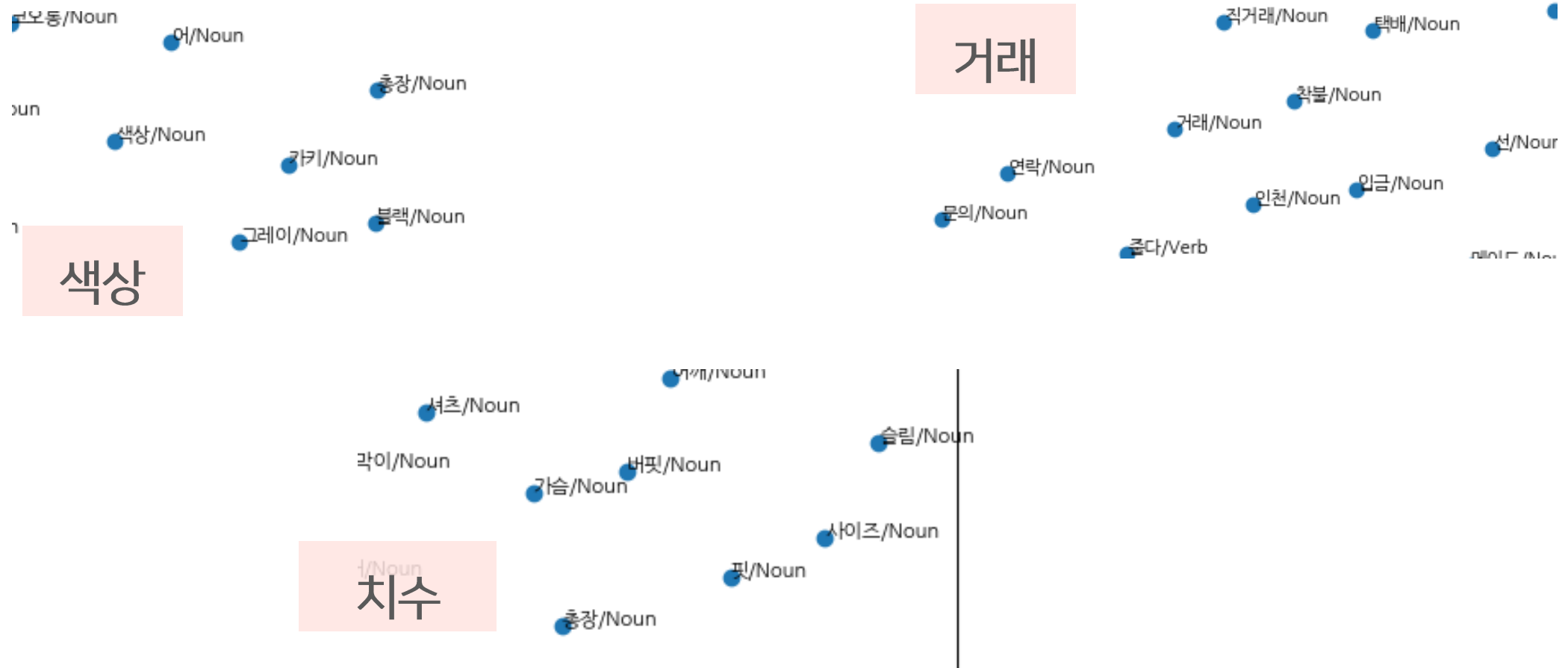


최소단어 수 / 윈도우 크기  
20 / 4

잘 출력됨

잘 출력됨

# 모델 결과





## 다음주 개발 계획

박시현

우분투 서버와 크롤링 프로그램 연결 설정  
통신 모듈 구현

박진영

웹 스토리지 구현  
서버 통신 구현 시작

이민희

서버 모듈 연결 / 모델 파라미터 변경  
상세 구조 정리

## 팀원 별 각 주간 개발 계획

박시현	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
주제 확정/제안서 작성								중간 고사 & 중간 데모 준비						최종 데모	최종 리포 트 제출
크롤링 · 전처리															
데이터셋 구축															
서버 구축															
통신 모듈 작업															
서버 통신 구현															
데모 준비 · 서비스 테스트															

## 팀원 별 각 주간 개발 계획

박진영	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
주제 확정/제안서작성								중간 고사 & 중간 데모 준비						최종 데모	최종 리포 트 제출
프론트엔드 화면 구현															
프론트엔드 동작 구현															
텍스트 대치 구현															
확장 앱 통신 구현															
웹 페이지 쿠키 구현															
테스트															

## 팀원 별 각 주간 개발 계획

이민희	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
주제 확정/제안서작성								중간 고사 & 중간 데모 준비						최종 데모	최종 리포 트 제출
자연어 처리 공부															
라이브러리 사용 전처리															
알고리즘 구체화															
문장분리															
주변 토큰 처리															
키워드 사전 변경															
모델 구축, 파라미터 변경															
서버 모듈 연결															
성능 및 정확도 향상															

공통: 연두색 / 박시현:핑크색 / 박진영:하늘색 / 이민희:주황색

[illegible]



감사합니다! QnA

