

# 감성 분석을 위한 토큰나이징 및 워드 임베딩 기법 탐색

20203065 소프트웨어학부 박규연



01

## Introduction

주제, 데이터 소개

02

## Preprocessing

데이터 전처리 및  
토큰화

03

## Tokenizer Comparison

형태소 분석기  
비교

---

# TABLE OF CONTENTS

## Word Embedding

Word2Vec와  
FastText 비교

04

## Keyword Analysis

워드 카운트  
워드 클라우드  
임베딩벡터 시각화

05

## Classification Model

감성 분류 모델  
훈련 및 성능비교

06



# 01

## Introduction

---

주제, 데이터 소개



# 감성 분석 모델 구축을 위한 토큰나이징, 워드 임베딩 기법 탐색



## 토큰화 기법 탐색

Kkma, Komoran, Okt, Mecab,  
KoNLTK

Word2Vec, FastText

## 워드 임베딩 기법 탐색



## 모델 훈련

Logistic Regression, LSTM

# NSMC Dataset

Naver sentiment movie corpus (20만 개)

 **nsmc** Public

Naver sentiment movie corpus

 Python  493  191



# 02

## Preprocessing

---

데이터 전처리 및 토큰화

# Preprocessing and Tokenizing

잼있어용~ㅋㅋㅋㅋㅋㅋㅋㅋ  
ㅋㅋ정일우♥♥♥♥♥♥♥♥



잼있어용ㅋㅋㅋㅋㅋㅋㅋㅋ  
ㅋㅋ정일우



정규 표현식을 이용한 특수문자 제거

`pattern=r'^a-zA-Z0-9가-힣 \n'`



명사 토큰화

`tagger.nouns(text)`



불용어 제거

`if token not in stopwords`

# Preprocessing and Tokenizing

잼있어용 ㅋㅋㅋㅋㅋㅋㅋㅋ  
ㅋ정일우



[ '잼', '있', '정일우' ]



정규 표현식을 이용한 특수문자 제거

`pattern=r'^a-zA-Z0-9가-힣 \n'`



명사 토큰화

`tagger.nouns(text)`

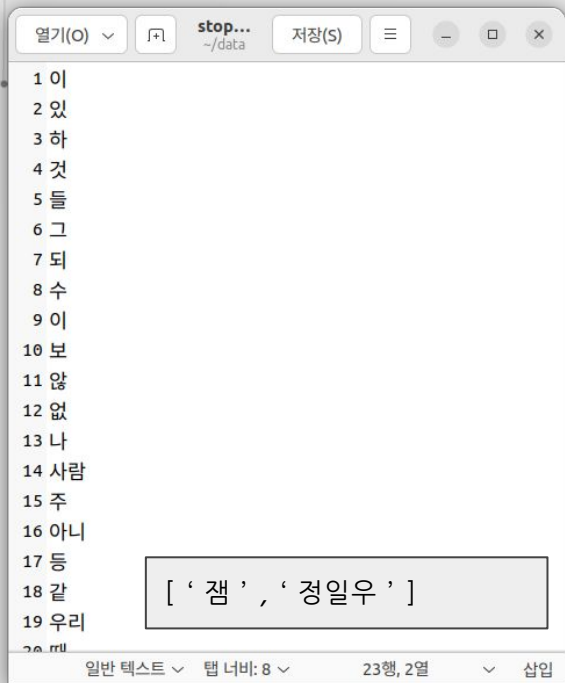


불용어 제거

`if token not in stopwords`



# Preprocessing and Tokenizing



정규 표현식을 이용한 특수문자 제거

`pattern=r'^[a-zA-Z0-9가-힣 \n]'`



명사 토큰화

`tagger.nouns(text)`



불용어 제거

`if token not in stopwords`

# Preprocessed Dataframe

	id	document	label	tokens
0	8112052	어릴때보고 지금다시봐도 재밌어요ㅋㅋ	1	['보고']
1	8132799	디자인을 배우는 학생으로, 외국디자이너와 그들이 일군 전통을 통해 발전해가는 문화산...	1	['디자인', '학생', '외국', '디자이너', '일군', '전통', '통해', ...]
2	4655635	폴리스스토리 시리즈는 1부터 뉴까지 버릴게 하나도 없음.. 최고.	1	['폴리스스토리', '시리즈', '부터', '뉴', '최고']
3	9251303	와.. 연기가 진짜 개쩔구나.. 지루할거라고 생각했는데 몰입해서 봤다.. 그래 이런...	1	['연기', '진짜', '몰입', '다그', '진짜', '영화']
4	10067386	안개 자욱한 밤하늘에 떠 있는 초승달 같은 영화.	1	['안개', '밤하늘', '초승달', '영화']
5	2190435	사랑을 해본사람이라면 처음부터 끝까지 웃을수 있는영화	1	['사랑', '라면', '처음', '끝', '영화']
6	9279041	완전 감동입니다 다시봐도 감동	1	['완전', '감동', '감동']
7	7865729	개들의 전쟁2 나오나요? 나오면 1빠로 보고 싶음	1	['전쟁', '빠']



# 03

## Tokenizer Comparison

---

형태소 분석기 비교

# Accuracy Comparison

	‘ 내 나이와 같은 영화를 지금 본 나는 감동적이다..하지만 훗날 다시보면대사하나하나 그 감정을완벽하게 이해할것만 같다... ’
Kkma	['내', '나이', '영화', '나', '감동', '훗날', '보면대', '사', '하나', '하나', '감정', '완벽', '이해', '것']
Komoran	['나이', '영화', '감동', '하지', '훗날', '다시', '보면', '대사', '하나하나', '감정', '을', '완벽', '이해', '것']
Okt	['내', '나이', '영화', '나', '감동적', '훗날', '대사', '대사하나하나', '하나하나', '감정', '완벽', '이해']
Mecab	['내', '나이', '영화', '지금', '나', '감동', '훗날', '다시', '사하나', '그', '감정', '이해']
KoNLTK	['나', '나이', '영화', '나', '감동적', '훗날', '다시보면대사하나', '감정을완벽', '이해할것']

Kkma, Komoran, Okt

띄어쓰기가 잘 되어있는 리뷰에서 괜찮은 성능을 보임

# Accuracy Comparison

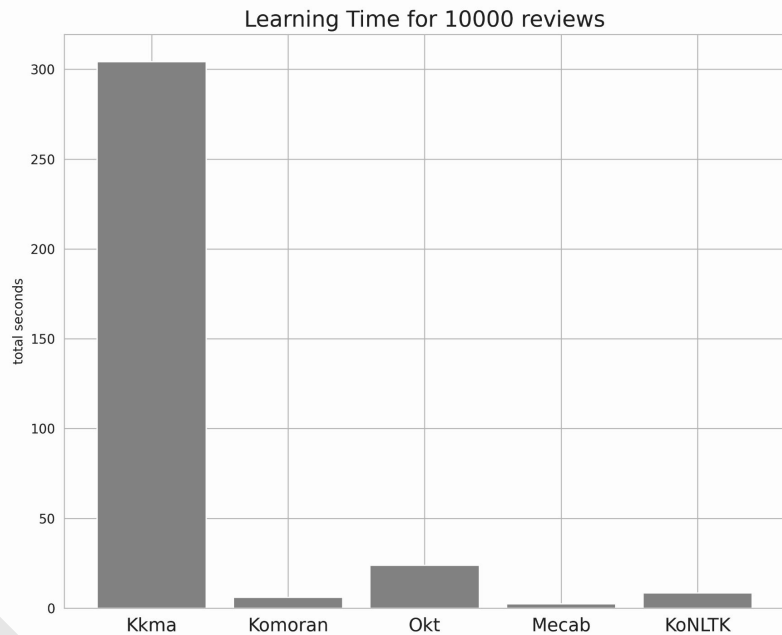
	'이거어렸을때되게재밋게봄 ㅋㅋ이정재 이범수 ㅋㅋㅋㅋ연기짬'
<b>Kkma</b>	['이거', '때', '이정재', '이범수', '연기']
<b>Komoran</b>	[]
<b>Okt</b>	['이거', '때', '봄', '이정재', '이', '이범수', '범수', '연기']
<b>Mecab</b>	['거', '때', '봄', '이정재', '이범수', '연기', '짬']
<b>KoNLTK</b>	['이거어렸을때되게재밋게봄', '이정', '이범수', '연기짬']

**Kkma**가 띄어쓰기가 잘 안 되어 있는 리뷰에서도 유의미한 토큰들을 잘 뽑아냄  
**Komoran**은 띄어쓰기가 잘 안되어 있는 리뷰에서는 성능이 떨어짐

# Speed Comparison

## Speed

**Mecab**이 속도 측면에서 매우 우수  
**Kkma**는 속도 측면에서 매우 떨어짐



# Tokenizer Comparison

	Kkma	Komoran	Okt	Mecab	KoNLTK
Accuracy	○	×	○	△	△
Speed	×	○	△	○	○

정확도 측면에서 가장 우수한 **Okt**  
속도 측면에서 가장 우수한 **Mecab**  
을 선정하여 이후 실험 진행

# 04

## Word Embedding

Word2Vec와 FastText 비교







# Most Similar Top 5

“ 배우 ”

	Word2Vec	FastText
Okt	‘ 연기자 ’ , ‘ 영화배우 ’ , ‘ 김혜수 ’ ‘ 양동근 ’ , ‘ 아역배우 ’	‘ 배우진 ’ , ‘ 단역배우 ’ , ‘ 재연배우 ’ ‘ 영화배우 ’ , ‘ 유명배우 ’
Mecab	‘ 연기자 ’ , ‘ 조연 ’ , ‘ 차승원 ’ ‘ 신인 ’ , ‘ 손예진 ’	‘ 파배우 ’ , ‘ 남배우 ’ , ‘ 명배우 ’ ‘ 영화배우 ’ , ‘ 여배우 ’



## Most Similar Top 5

positive=[ “ 감동 ” , “ 마지막 ” ]

	Word2Vec	FastText
Okt	‘ 핑 ’ , ‘ 콧 ’ , ‘ 피아노 ’ ‘ 글썽 ’ , ‘ 주룩주룩 ’	‘ 마지막모습 ’ , ‘ 마지막방송 ’ , ‘ 핑 ’ ‘ 글썽 ’ , ‘ 방울 ’
Mecab	‘ 콧물 ’ , ‘ 방울 ’ , ‘ 흘 ’ ‘ 물결 ’ , ‘ 하이라이트 ’	‘ 흘 ’ , ‘ 콧물 ’ , ‘ 방울 ’ ‘ 감동이 ’ , ‘ 핑 ’



## Most Similar Top 5

positive=[ “ 황정민 ” , “ 최우식 ” ]

negative=[ “ 별로 ” ]

	Word2Vec	FastText
Okt	‘ 김영호 ’ , ‘ 경력 ’ , ‘ 김갑수 ’ ‘ 백성현 ’ , ‘ 이성민 ’	‘ 상우 ’ , ‘ 백성현 ’ , ‘ 손현주 ’ ‘ 김영호 ’ , ‘ 정려원 ’
Mecab	‘ 차승원 ’ , ‘ 디카프리오 ’ , ‘ 엄정화 ’ ‘ 손예진 ’ , ‘ 임창정 ’	‘ 김혜수 ’ , ‘ 차승원 ’ , ‘ 안성기 ’ ‘ 엄정화 ’ , ‘ 신인 ’

The top corners of the slide feature decorative geometric patterns. On the left, there are several interconnected triangles and lines, some with small black dots at their vertices. On the right, a more complex network of lines and dots is visible, resembling a graph or a molecular structure.

# 05

## Keyword Analysis

워드카운트 / 워드클라우드 / 임베딩벡터 시각화

# Word Count

`count.most_common(8)`

**Okt**

[('영화', 57883), ('정말', 11748), ('진짜', 10643), ('평점', 8394), ('최고', 7805), ('연기', 7047), ('스토리', 6734), ('이영화', 6516)]

**Mecab**

[('영화', 75236), ('연기', 8901), ('최고', 8586), ('평점', 8238), ('스토리', 7020), ('드라마', 6703), ('감동', 6309), ('배우', 5643)]

\*Okt 사용

# Word Cloud

가장 스토리 추천 긴장감 최고 대한 추억  
감동 재미 인간 밈 제멋대로 보고  
표현무엇 조금엔덜 노래 사랑 가족 기분 처음 드라마 기도  
마지막 연기 액션 평점 만점 이야기 잘만 아이  
주인공 정말 최고다 그냥이제 결말 인생 파  
진짜 연극 공감 한국 내 물인 영화 계속 현실  
작품 보기 이영화 배 배우 기대애디  
한번 역시 여운 기억 모습 장면 이해 마음 매력 반전

Positive

Negative

하반기 연극 시리즈 감동 감독 자체 처음 작가 진짜 이견  
연출 힘있 공도영화 출연 시작 열주 열왕 절대 연기력 현실합력타원  
진심 관객 데 영화 스토리 시작 열왕 열왕 열왕 열왕  
정말 나 대가 아이 노서 열왕 열왕 열왕 열왕  
소재 차라리 이해제목 액션 전도 열왕 열왕 열왕 열왕  
원작 보고 용도 오빠 완전 기도 평점 못 데 배 영화  
스토리아 최고 완전 기도 평점 못 데 배 영화  
이야기 수준 캐릭터 코미디 설정 노점 쓰레기 등

# Word Cloud



Positive

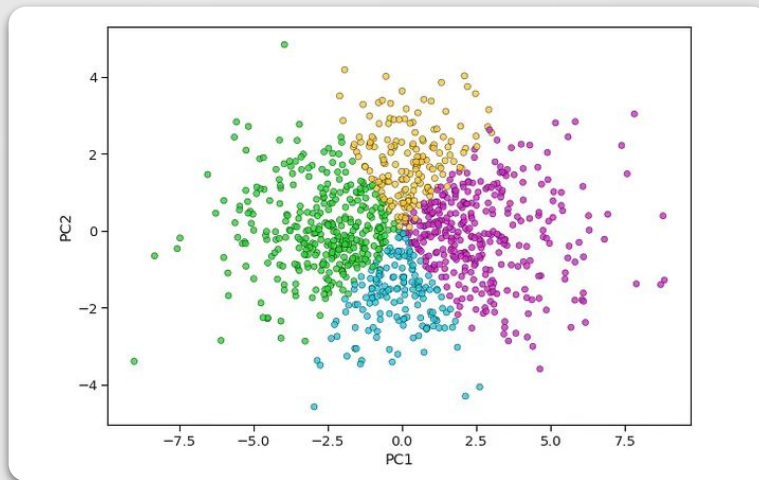
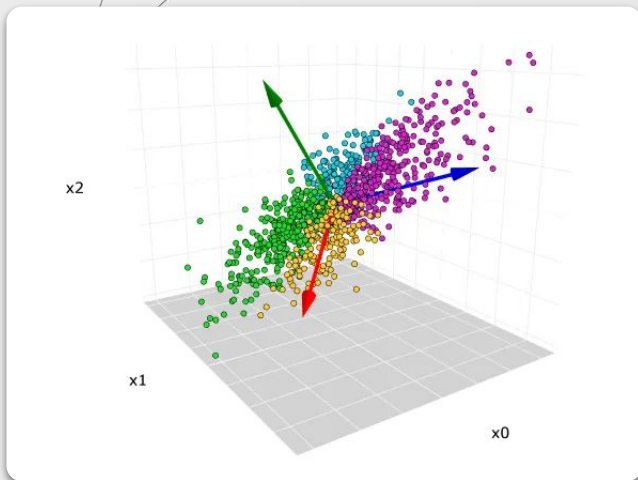
Negative



지배적인 단어,  
키워드 분석에 의미 없는 단어 제외

# PCA (주성분 분석)

고차원의 데이터를 저차원의 데이터로 환원시키는 기법

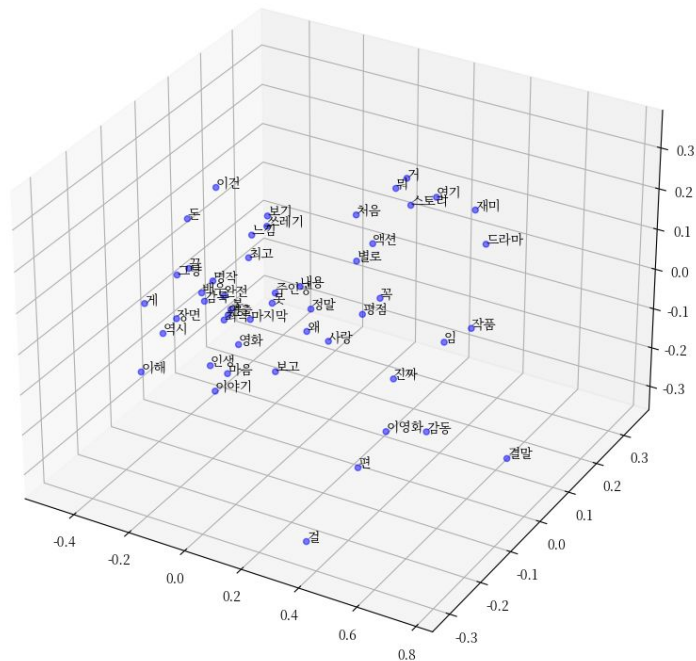
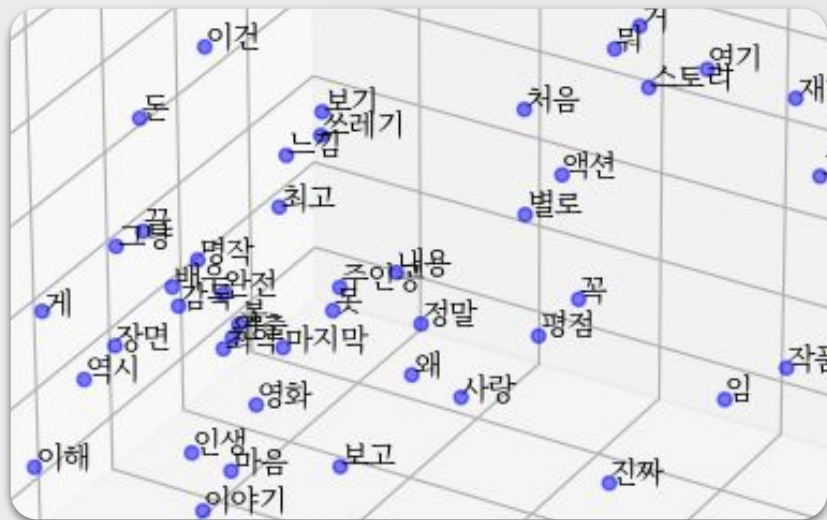




# 3D Visualization

## Most Similar Top 100

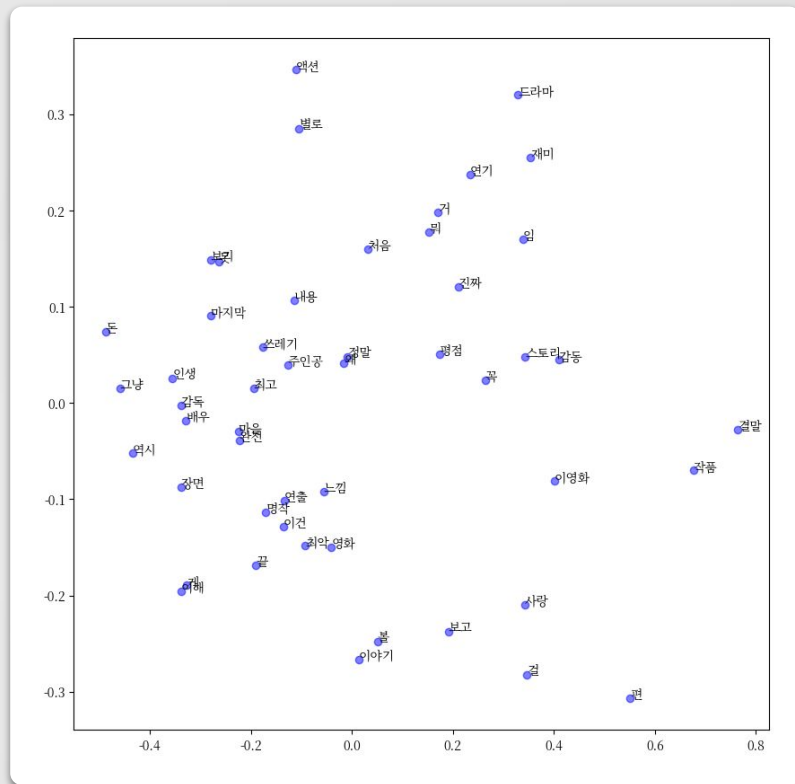
positive=[ “ 황정민 ” , “ 마동석 ” ]



# 2D Visualization

## Most Similar Top 100

positive=[ “ 황정민 ” , “ 마동석 ” ]





# 06

## Classification Model

감성 분류 모델 훈련 및 비교

# Modeling

The slide features decorative geometric patterns. The top-left corner has a network of interconnected lines and dots, with several triangles of varying sizes. The bottom-right corner contains a cluster of small, scattered circles.

**Models**

A diagram illustrating the relationship between a general category and specific models. A central light gray box labeled 'Models' has two lines extending downwards from its bottom edge. Each line connects to a dark gray box. The left box is labeled 'Logistic Regression' and the right box is labeled 'LSTM'.

```
graph TD; Models[Models] --- LR[Logistic Regression]; Models --- LSTM[LSTM]
```

**Logistic  
Regression**

**LSTM**



# Embedding

**Word2Vec  
(CBOW)**

```
[[1, 1, 1, 1, 1, 1, 1, 1, 1, 1],  
 [1, 2, 1, 3, 1, 1, 1, 1, 1, 1],  
   ...  
 [0, 1, 2, 1, 3, 2, 1, 3, 2, 1]]
```

**Word2Vec  
(Skip-gram)**



```
[1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
```

**FastText**

단어 벡터들을 평균내서 문장 벡터 생성

# Logistic Regression

	Word2Vec (CBOW)	Word2Vec (Skip-gram)	FastText
Okt	69.61%	72.00%	69.45%
Mecab	69.48%	71.00%	69.04%

# LSTM

	Word2Vec (CBOW)	Word2Vec (Skip-gram)	FastText
Okt	63.28%	72.67%	70.45%
Mecab	69.58%	71.08%	70.00%

# Conclusion

**Okt**

Tokenizer

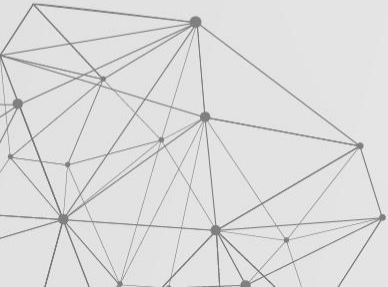
**Word2Vec**

(Skip-gram)

Word Embedding

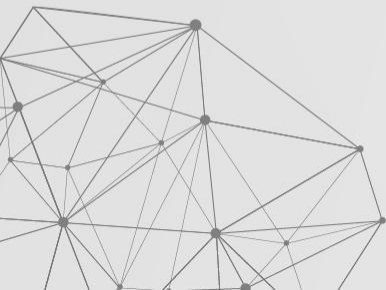
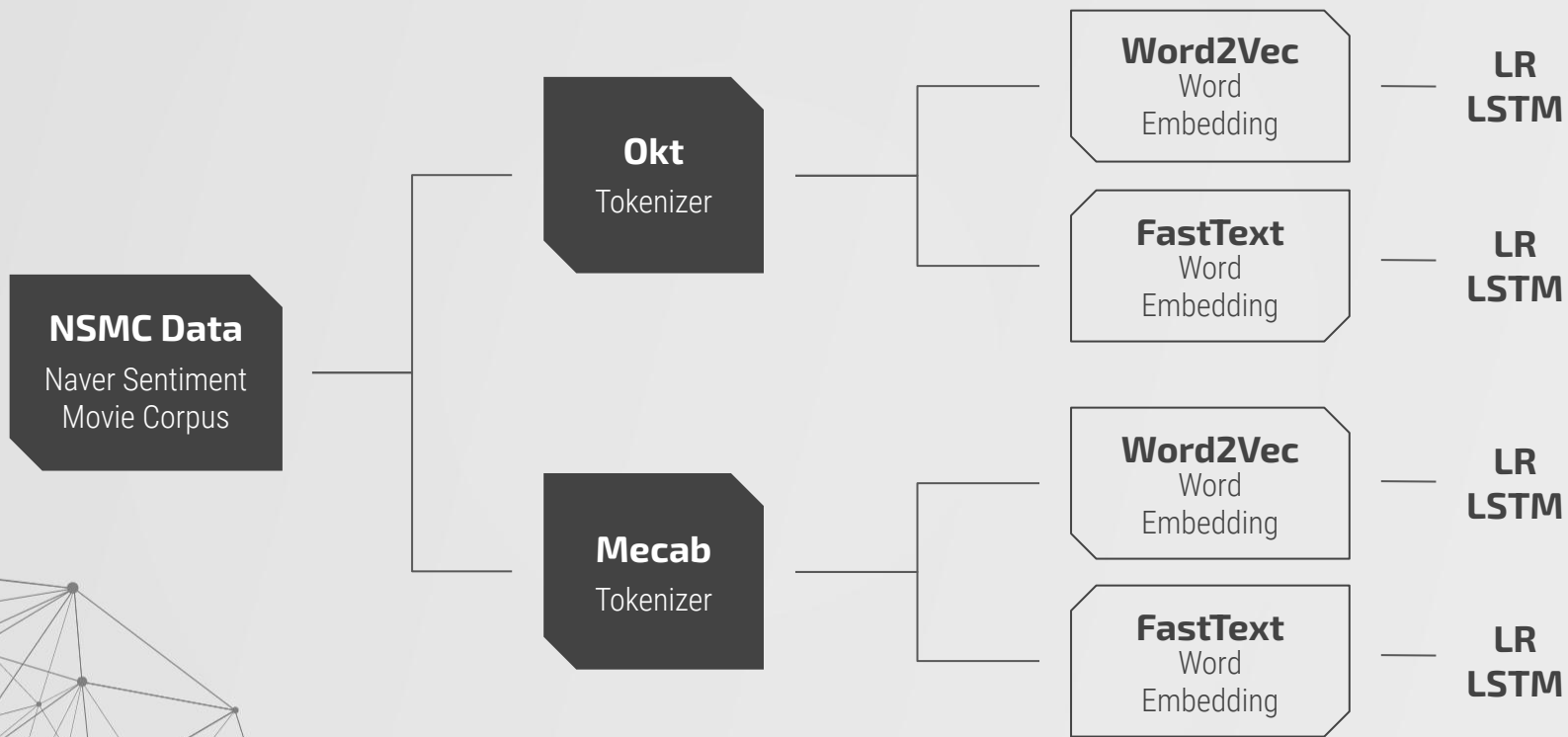
**LSTM**

Classification Model





# Process



# References

- ChatGPT
- Word2vec 및 fastText 임베딩 모델의 성능 비교, 2020, 강현석 외 1명, 디지털콘텐츠학회논문지
- 감성 분류를 위한 워드 임베딩 성능 비교, 2021, 윤혜진 외 2명, ACK 2021 학술발표대회 논문집





<https://github.com/noooey/Exploring-for-Sentiment-Analysis>

# THANKS

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.

**Please keep this slide for attribution.**