

# Final Project 기획안

5조

Leader 최호진  
김한준  
김선림  
이재원  
윤선영  
최가은



# Intro

---

## 1. 주제선정 및 의의

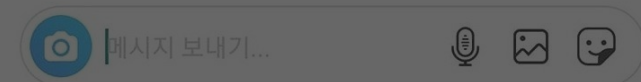
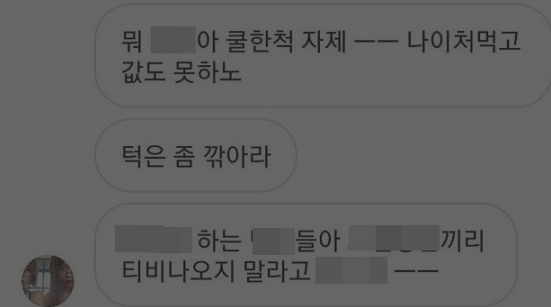
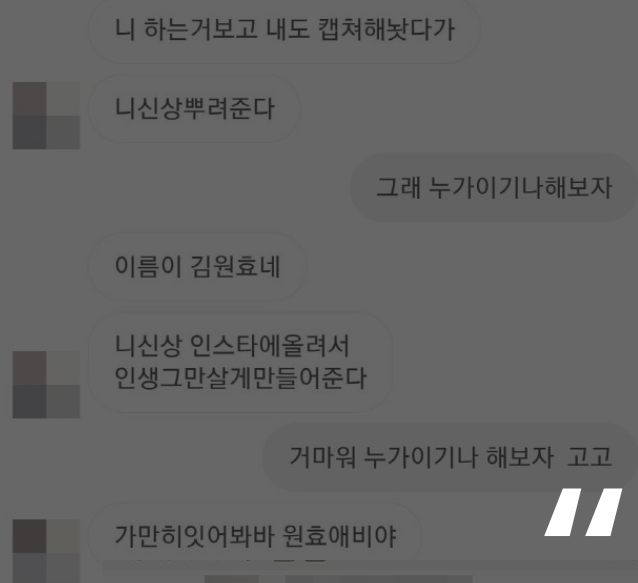
## 2. 분석계획

- 크롤링
- 악성댓글 여부 라벨링
- 텍스트 전처리
- 머신러닝 & 딥러닝
- Django Web Framework

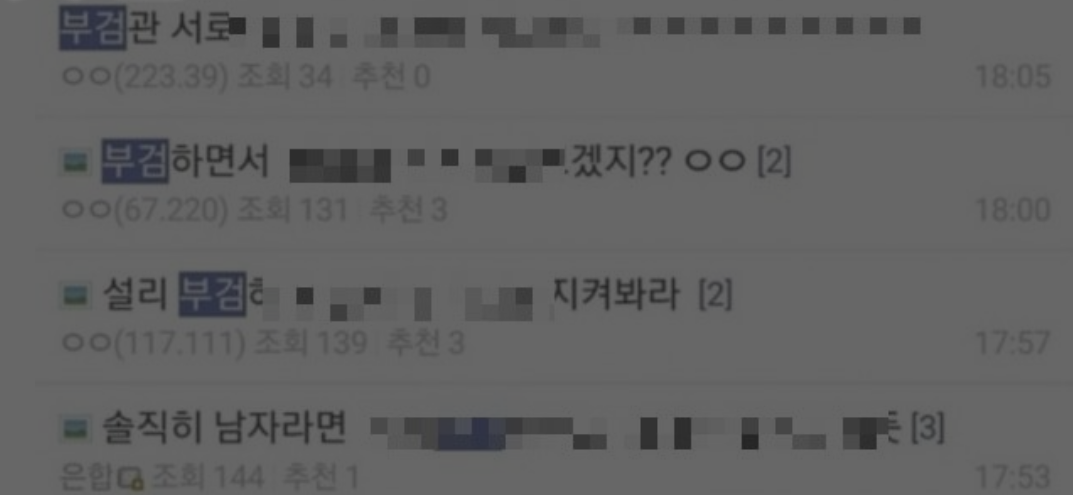
# 1. 주제선정 및 의의

---

## 1. 주제선정 및 의의



## 악성 댓글



## 1. 주제선정 및 의의

[기자수첩] 네이버 연예뉴스 댓글이 폐지됐다 - 민중의소리

네이버 스포츠뉴스 댓글 잠정 중단... "악성

포털 다음, 연예 뉴스 댓글 잠정 폐지..."인격 모독

연예계 "공포의 댓글창 폐지 환영..."

## 1. 주제선정 및 의의

연예인 이어 일반인도 노린다... '악플민국'의 민낯

압력 2020.11.08 17:09 | 수정 2020.11.09 00:40 | 자면 A33

사회 >

동문들로부터 악성 댓글받은 대학생 극단적 선택...경찰  
“고소장 접수”

소리없는 흥기 '악플' 공세...연예인 일반인 안가린다

“글 하나 올렸는데 악플 1000개” 죽음 부르는 ‘사이버불링’

//

# 악플은 해결이 시급한 사회 문제

//

연예인 이어 일반인도 노린다... '악플민국'의 민낯

입력 2020.11.08 17:09 | 수정 2020.11.09 00:40 | 자면 A33

사회 >

동문들로부터 악성 댓글받은 대학생 극단적 선택...경찰

소리없는 흥기 '악플' 공세...연예인 일반인 안가린다

“글 하나 올렸는데 악플 1000개” 죽음 부르는 ‘사이버불링’

## 1. 주제 선정

//

머신러닝 혹은 딥러닝을 이용한 악성댓글 탐지 봇

//

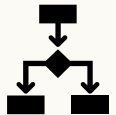


## 1. 주제선정 및 의의

### 지난 4개월간 배운 것



Crawling



Machine Learning + Deep Learning



Django Web Framework

### “악플 감지 봇” 프로젝트 적용

Crawling을 통한 댓글데이터 수집

악플 분류

서비스를 위해 Web에 구현

## 2. 분석계획

---

### - 크롤링

:댓글 데이터를 수집하기 위해 '네이버' or '인스타그램'의 댓글을 크롤링



BeautifulSoup

### - 악성댓글 여부 라벨링

INDEX	댓글 내용
1	당신의 앞날을 응원합니다.
2	너무 멋있어요 ㅎㅎ
3	야, 이 <b>미친 새끼야</b> , OO가 아니고 XX 이게 맞아. 제대로 알고 떠들어.
4	'OO이 아니라 XX이기에 수정했습니다. 출처가 불분명한 정보는 서술하지 말아주세요.'

객관적인 기준으로 악플라벨링



악플여부
0
0
1
0

### -텍스트 전처리 텍스트데이터 토큰화

: 크롤링 한 댓글을 '단어' or '형태소' or '음절' 형태로 토큰화를 진행하기 위한 작업

#### 토큰화 예시

[열심히 일한 당신, 연휴에는 여행을 가봐요]



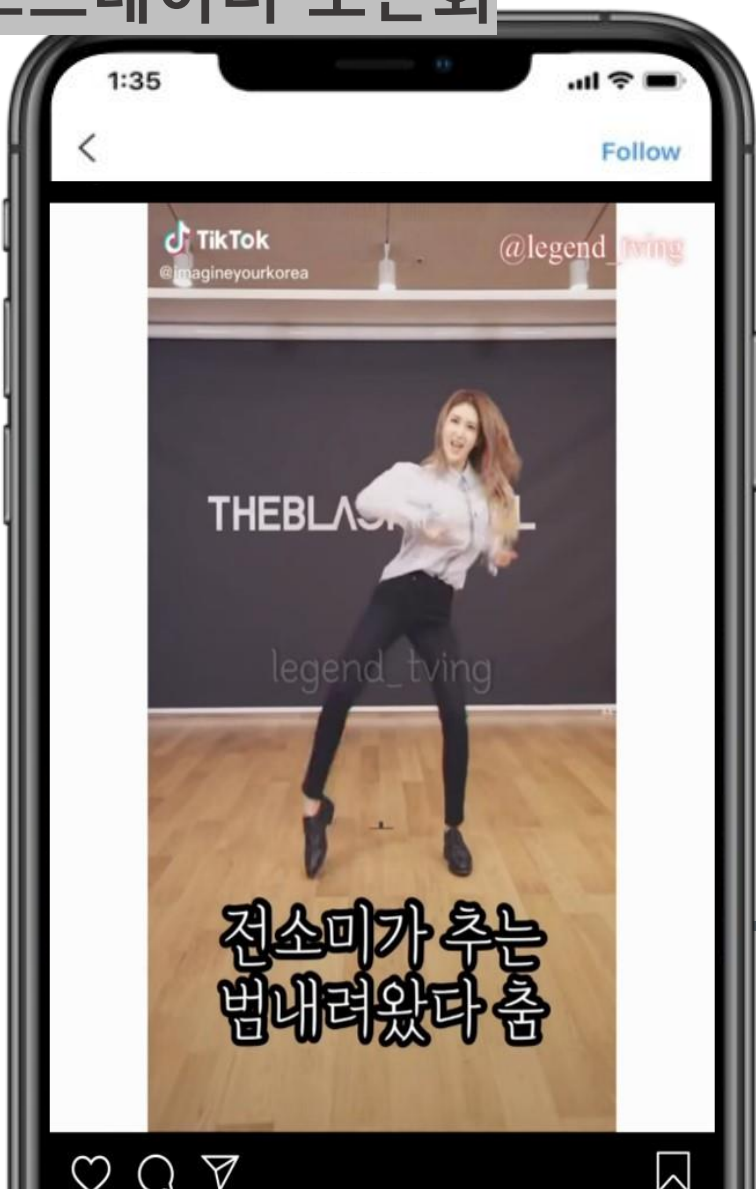
형태소 단위의 토큰화: ['열심히', '일한', '당신', '연휴', '에', '는', '여행', '을', '가보', '아요']

단어 단위의 토큰화: ['일', '당신', '연휴', '여행']

-텍스트 전처리    텍스트데이터 토큰화

:댓글 데이터 특징

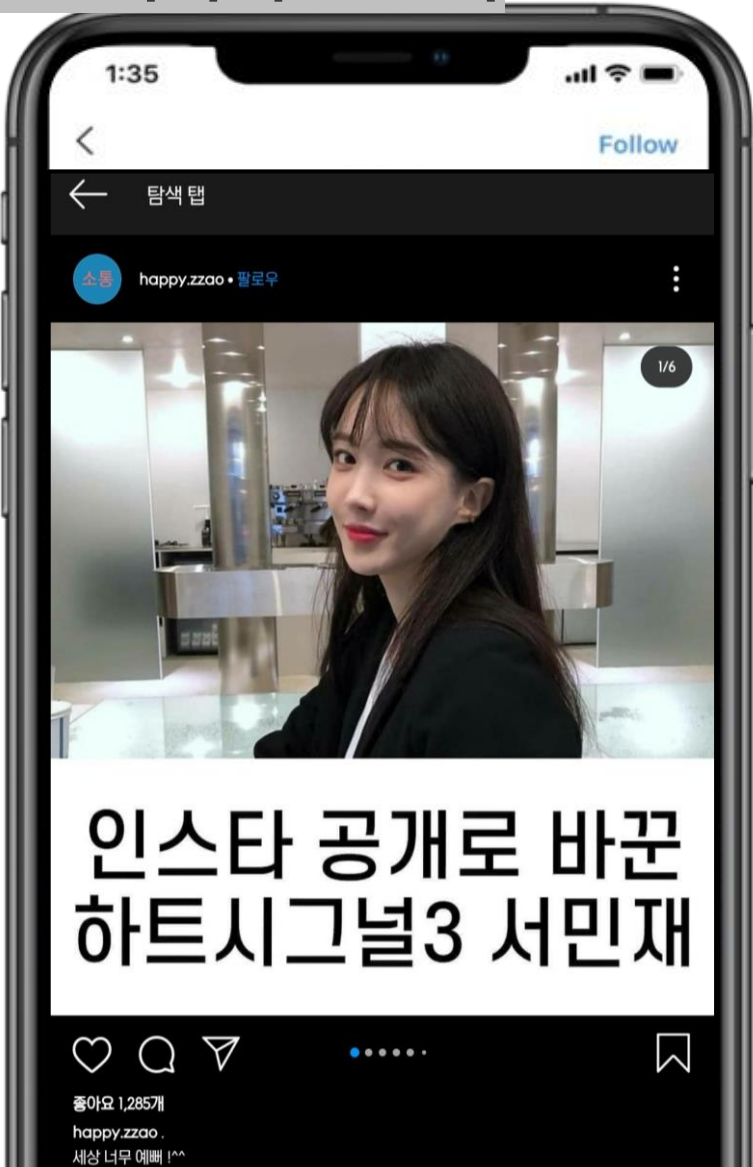
1. 간결함



-텍스트 전처리    텍스트데이터 토큰화

:댓글 데이터 특징

2.    맞춤법 X



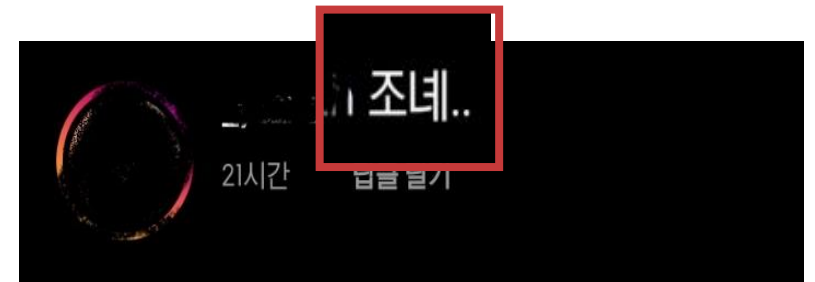
### - 텍스트 전처리 텍스트데이터 토큰화

: 댓글 길이가 짧고 맞춤법이 고르지 못할 때는 '단어', '형태소' 단위로 토큰화 진행이 어려움

So, 음절단위 토큰화 필요



'범', '내', '려', '온', '당'



'조', '네'



### - 텍스트 전처리 텍스트데이터 토큰화

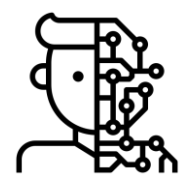
// 만약, 음절단위의 토큰이 좋은 성능을 보장하지 않는다면? //

Twitter에서 제공하는 Okt , 카카오톡에서 제공하는 Khaiii 등의  
SNS 기반 형태소 분석기 사용



-텍스트 전처리 토큰 Labeling

: 만들어진 토큰을 머신러닝 또는 딥러닝에 적용하기 위해 고유 인덱스 번호를 부여



토큰 Labeling

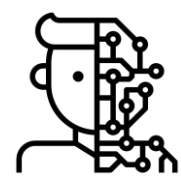
토큰
'범', '내', '려', '온', '다'
'조', '네'



토큰 라벨링
'1', '2', '3', '4', '5'
'6', '7'

- 텍스트 전처리    토큰 Padding

: 토큰화와 라벨링이 이루어진 토큰을 의미없는 숫자(0)을 이용하여 데이터의 길이를 동일하게 만드는 과정



토큰 Padding : 텍스트데이터를 DataFrame화 하는 과정

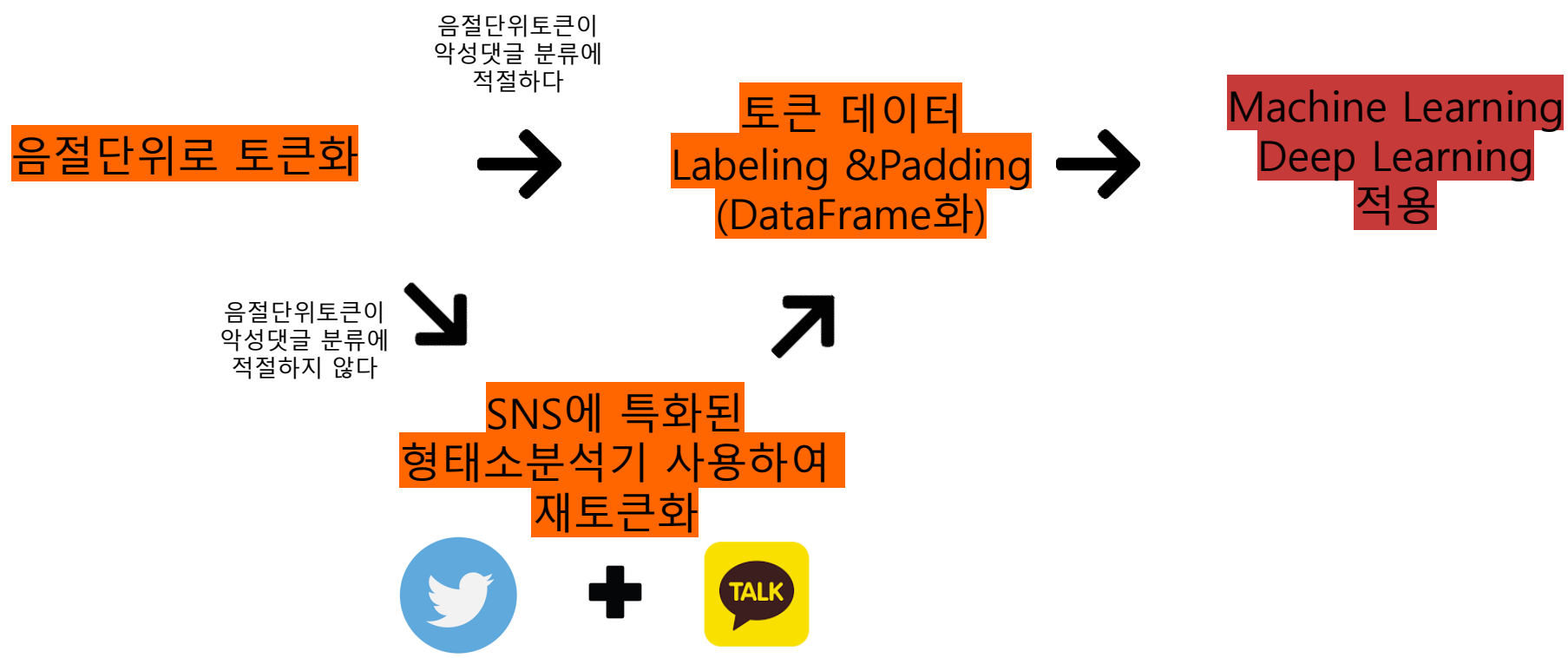
토큰 라벨링
'1', '2', '3', '4', '5'
'6', '7'



Index	음절1	음절2	음절3	음절4	음절5
row1	1	2	3	4	5
row2	6	7	0	0	0

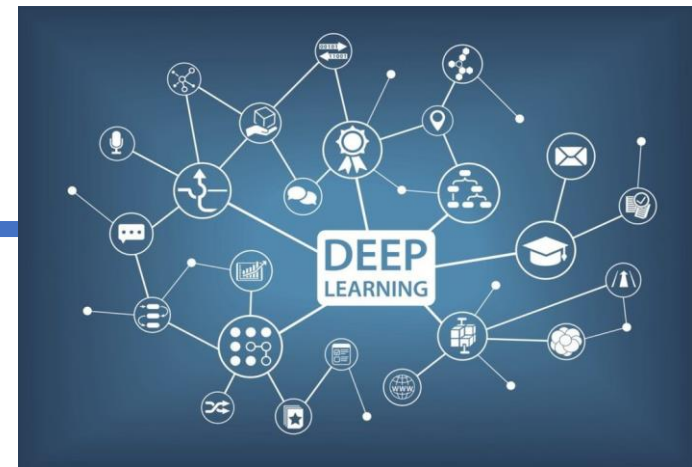
- 텍스트 전처리 한눈에보기

한눈에 보는 텍스트전처리



### -악성댓글분류 머신러닝 & 딥러닝

최근 동향으로 자연어처리(NLP)는 머신러닝 보다 딥러닝이 더욱 많이 사용됨

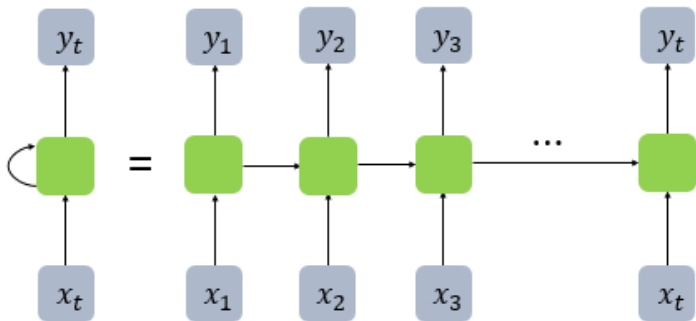


-악성댓글분류 머신러닝 & 딥러닝

특히, Google에서 개발한 Bert라는 딥러닝 모델의 성능이 우수



만약 딥러닝을 사용한다면 RNN계열(RNN, LSTM) 모델을 사용할 예정



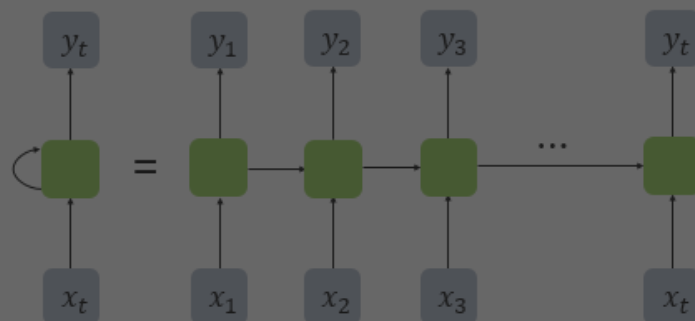
### -악성댓글분류 머신러닝 & 딥러닝

특히, Google에서 개발한 Bert라는 딥러닝 모델의 성능이 우수

“ BUT 딥러닝은 머신러닝에 비해  
비교적 많은 양의 데이터를  
필요로 함 ”



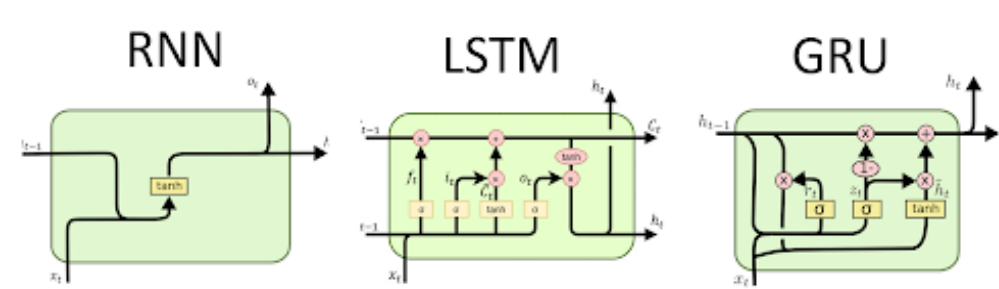
만약 딥러닝을 사용한다면 RNN (RNN, LSTM) 모델을 사용할 예정



-악성댓글분류 머신러닝 & 딥러닝

딥러닝 시나리오

- 1. 프로젝트 기간 동안 많은 데이터를 크롤링
- 2. 많은 데이터를 이용해 RNN 계열(LSTM) 딥러닝을 주 모델로 사용하여 악성 댓글 여부를 예측



머신러닝 시나리오

- 1. 프로젝트 기간 동안 많은 데이터를 크롤링
- 2. 데이터를 이용해 여러 머신러닝 기법에 적용하여 예측





## -Django WebFrame

악플 No



← →

A group photo of seven people, four women and three men, posing together. The man in the center is wearing a white shirt with a tiger pattern and a face mask. The man on the left is wearing a black shirt with a tiger pattern. The man on the right is wearing a black shirt with a tiger pattern. The women are wearing various clothing, including a red patterned top, a black top, and a light blue denim jacket.

사진 너무 잘 나왔다아~

입력

## -Django WebFrame

악플 No



## -Django WebFrame

악플 Yes



## -Django WebFrame

악플 Yes



THANK YOU

악! 플원정대



# THANK YOU

## 악!픈원정대

produced by 지미 호

