

49

멀티모달 기반 스팸 필터링 플랫폼 개발

소속 정보컴퓨터공학부

분과 D

팀명 멀티모달

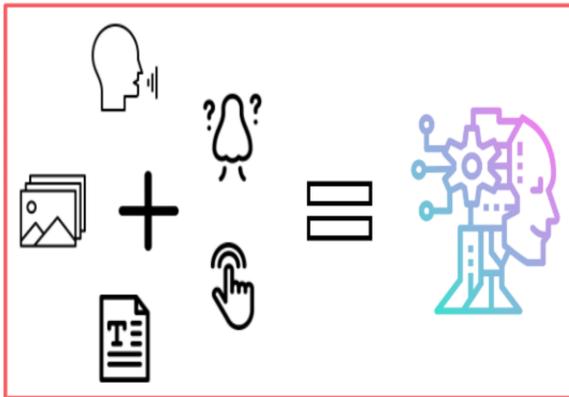
참여학생 이강우, 윤상호, 조재홍

지도교수 최윤호

과제 소개 및 목표

과제 소개 - 멀티모달이란?

- **모달리티** 라는 것은 데이터의 형태, 유형을 의미한다. 예를 들어, 언어라는 것은 시각적으로 글로 표현될 수 있고, 청각적으로 소리로 표현될 수 있다. 또한 **점자는 촉각의 형태로 언어를 표현한 것이다.** 이처럼 시각, 촉각, 청각 등을 하나의 모달리티라고 할 수 있다.
- 앞서 소개한 기사에 따르면 요즘은 대부분의 텍스트 필터링을 우회하기 위해서 이미지 형식으로 스팸메일을 보내는 경우가 많아지고 있다.
- 본 졸업과제에서는 멀티모달을 활용하여 수집한 텍스트 메일뿐만 아니라, 이미지와 Chat-GPT로 생성한 텍스트 메일을 각각의 모달리티로 보고 이 3가지 모달을 잘 처리할 수 있는 스팸 필터링 모델을 개발하였다.



결론은, 텍스트와 이미지가 섞여 들어와도 잘 처리할 수 있다!

- 최근 텍스트 스팸은 대부분의 필터링 시스템에서 잘 처리되고 있다.
- BUT, 최근에는 텍스트 필터링을 우회하기 위해서 **이미지 형식으로 보낸다거나, 이미지와 텍스트를 같이 동봉하여 전송해서 기존의 필터링 시스템을 우회하는 스팸머들이 많이 등장하고 있다.**

'문자' 걸러내니 '이미지'·'카톡'으로...스팸의 진화



과제 목표

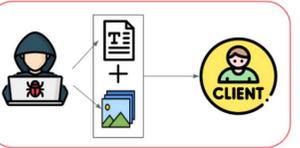
- 멀티모달리티를 활용하여 멀티미디어 형식으로 전송되는 스팸메일 필터링을 목표로 한다.
- 각 모달리티에 대한 피처를 직접 분석 및 선정해 새로운 유형의 스팸메일에 유연하게 대처할 수 있다.
- 필터링 결과를 눈으로 확인할 수 있도록 시각화 인터페이스 (스팸메일 여부, 스팸 분류 이유 등 유의미한 정보)를 구현한다.

과제 내용

동작 과정

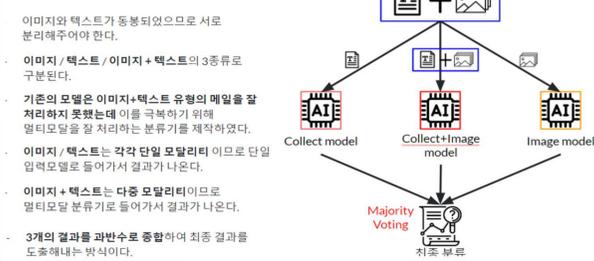
동작 과정 - 이미지 + 텍스트 입력

- <주어진 상황>
- 어떤 스팸머가 악의적인 의도를 가지고 스팸 필터링을 우회하려고 한다.
- 실제로 네이버와 구글에 정상적인 텍스트와 스팸 이미지를 같이 동봉해서 전송했을 경우 자동 필터링이 되지 않는 것을 확인했다.



- 본 졸업과제에서 사용하는 5가지 모델(분류기)
- 단일 입력 텍스트 메일 분류기
- 단일 입력 이미지 메일 분류기
- 멀티모달 텍스트+이미지 메일 분류기
- 멀티모달 생성+이미지 메일 분류기

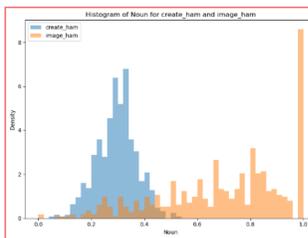
동작 과정 - 이미지 + 텍스트 분류



분석 및 결과

분석 - 피처 분석 방법

- 피처들의 각 모달리티간의 차별적 여부를 시각적으로 확인하기 위해 히스토그램을 그려 분포를 확인
- 각 모달리티의 스팸메일 간의 피처 분포도를 보면 답이 나오지 않을까? 라는 생각을 했고 실제로 히스토그램으로 그려서 살펴보면, 두 개의 그래프가 겹치는 면적이 50% 미만인 것을 선택하면 두 모달리티간의 차이를 명확하게 표현할 수 있다고 판단하였다.

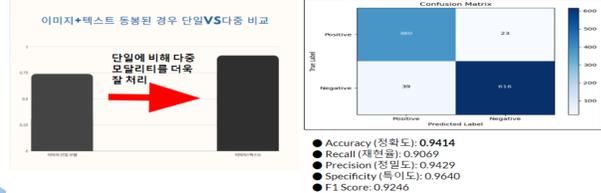


	생성+이미지(text)	생성+이미지(image)
special_ratio	0.527458768	0.493104871
number_ratio	0.320856888	0.355738076
self_count	0.46703034	0.51841721
upper_ratio	0.1480129	0.130628516
blank_ratio	0.772814951	0.306008094
crf_ratio	0.400054735	0.284124516
Noun	0.579741564	0.179713874
Pronoun	0.316468297	0.172116825
Verb	0.705450995	0.335550023
Adjective	0.461977787	0.389804043
Adverb	0.644791817	0.305467022
avg_word_sentences	0.494451361	0.549831774
avg_char_sentences	0.471555175	0.523792021
avg_word_paragraphs	0.111520348	0.548824281
avg_char_paragraphs	0.106599339	0.556219024

각 feature별 겹치는 비율을 나타낸 표

결과 - 최종 성능(이미지)

아래 그래프는 이미지 단일 모델과, 이미지 다중 모델(다중 모달리티)에 이미지 + 텍스트를 넣었을 때 성능을 나타낸 Confusion Matrix입니다.



WEB, DATABASE

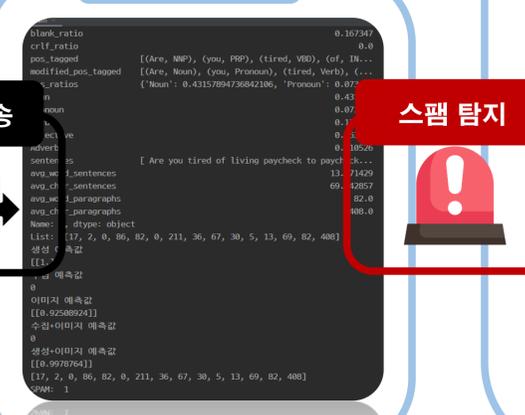


과제 결과

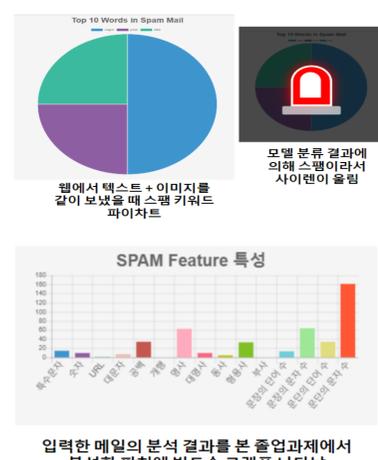
멀티모달 메일



메시지 분석



분석 결과



분석 결과

웹에서 텍스트 + 이미지를 같이 보냈을 때 스팸 키워드 파이차트

모델 분류 결과에 의해 스팸이러사 사이렌이 울림

입력한 메일의 분석 결과를 본 졸업과제에서 분석한 피처에 빈도수 그래프 나타남