

목차 A table of contents

- **1** 주제 소개
- 2 진행상황
- 3 향후계획



Part 1 주제 소개

가짜 뉴스란?

ADOUL 307,000 TESUILS (0.37 SECORUS)

https://www.donga.com > news > Society > article > all •

마늘-카레 먹으면 예방? 코로나 뺨치는 '인포데믹': 뉴스: 동아닷컴

Mar 18, 2020 — 실제 온라인에는 **코로나**19를 둘러싸고 확인되지 않은 정보가 난무하고 있다. 소금물을 비롯해 알코올, 마늘, **카레**를 섭취하면 바이러스를 **예방할** 수 ...

https://news.joins.com → article ▼

14억 명 중 확진자 3명, 인도가 코로나19 강한 건 카레 덕분 ...

Mar 3, 2020 — 그렇다면 인도 사람들은 정말로 **코로나**에 걸리지 않았을 - 신종 **코로나**바이러스, 인도사람,**코로나**,**코로나** 확진자,인도 전역,강황가루,인도식 **카레**, ...

About 1,510,000 results (0.40 seconds)

https://m.health.chosun.com > svc > news_view •

"한국의 김치, 코로나 사망률 낮춰"... 바이러스 침투 막는다

Jul 20, 2020 — 한국이 다른 나라에 비해 신종 **코로나**바이러스 감염증(**코로나**19) 사망자 수가 적은 이유가 '**김치'** 때문이라는 분석이 나왔다.

https://www.kita.net > cmmrcNews > cmmrcNewsDetail •

러 언론 "김치, 韓 코로나 발생·사망률 낮춰"...수출에 날개달까 ...

Apr 9, 2021 — 한국의 **김치**가 **코로나**19 **예방**과 증상 완화에 기여한다는 기사가 최근 러시아에서 대대적으로 보도됐다. **김치** 유산균의 항**코로나**바이러스 효능을 ...

http://www.kfdn.co.kr > ... ▼

[식약일보] 김치 코로나19에 미치는 효과...러시아 국영통신 최근 ...

Apr 6, 2021 — 중국이 뜬금없이 **김치** 원조설을 주장하고 있는 가운데 한국의 **김치**가 **코로나**19 **예** 방과 증상 완화에 이바지한다는 기사가 최근 러시아에서 대대적으로 ...

문제 배경

정보의 중요성 증가

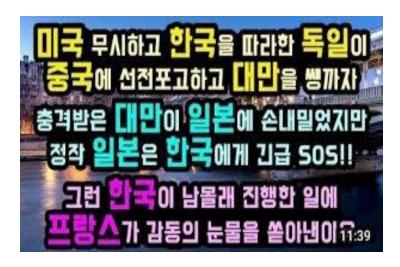
21세기 정보화 시대

인터넷 보급 증가

정보에 손쉬운 접근이 가능해짐



가짜 뉴스 사례 - 유튜브







문제점

- 전문성이 결여된 가짜 뉴스가 전세계적으로 많이 생산되고 있음
- 이런 가짜 뉴스를 진지하게 믿는 사람들이 많이 나타남

가짜 뉴스와 관련된 연구

Original research

YouTube as a source of information on COVID-19: a pandemic of misinformation? 8



b Heidi Oi-Yee Li¹, b Adrian Bailey¹, b David Huynh^{2, 3}, James Chan⁴ Correspondence to Ms Heidi Oi-Yee Li; heidi.li@live.ca; Ms Heidi Oi-Yee Li; heidi.li@live.ca

Abstract

Introduction The COVID-19 pandemic is this century's largest public health emergency and its successful management relies on the effective dissemination of factual information. As a social media platform with billions of daily views, YouTube has tremendous potential to both support and hinder public health efforts. However, the usefulness and accuracy of most viewed YouTube videos on COVID-19 have not been a like of the investigated.



Log in Q =

BRIEF REPORT | VOLUME 38, ISSUE 3, E1-E3, MARCH 01, 2010

YouTube As a Source of Information on the H1N1 Influenza Pandemic

Ambarish Pandey, MBBS 🖇 🖾 • Nivedita Patni, MBBS • Mansher Singh, MBBS •

Akshay Sood, MBBS • Gayatri Singh DOI: https://doi.org/10.1016/j.amepre.2009.11.007



연구 결과

- 코로나19를 키워드로 한 영상 중 조회수가 높은 <mark>69개의 영상 중 19개(약 27%)</mark>가 완전한 가짜 뉴스를 전달하고 있었음.
- H1N1(인플루엔자)가 유행한 시기에도 <mark>23%</mark>의 유튜브 영상들이 가짜 뉴스를 전달하고 있었음.

프로젝트 목표 + 기대효과

 STEP 1
 STEP 2
 STEP 3

 가짜 뉴스 분류기 구축
 >> 보류기 학습
 >> 통한 객관적이고 올바른 정보 제공

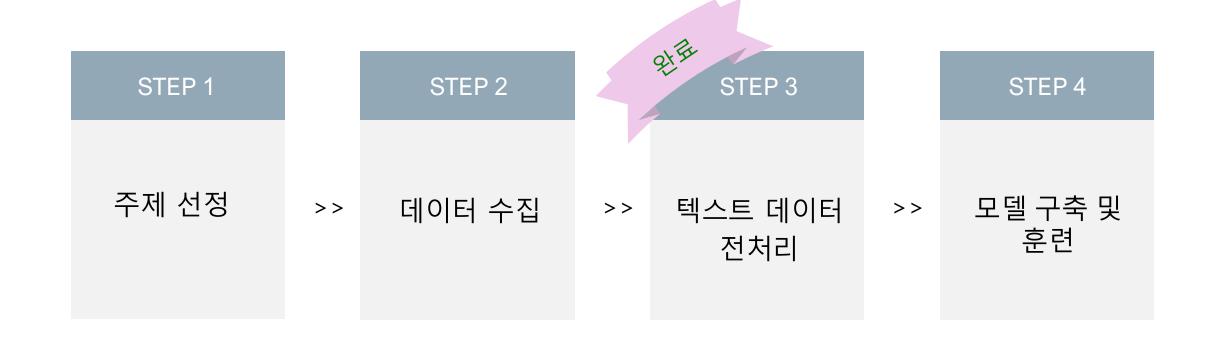
기대효과

- 가짜 뉴스 분류기를 구축함으로써 가짜 뉴스와 진짜 뉴스를 손쉽게 구분하고 가짜 뉴스로 인해 생기는 피해를 줄일 수 있다
- 가짜 뉴스를 구별하기 힘들어하는 사람들에게 도움을 줄 수 있다

 Part 2

 진행 상황

진행 상황



데이터설명

- 데이터 종류: 텍스트 데이터
- 데이터 구성 : title(제목), text(내용), subject(주제), date(날짜)
 - 데이터 속 가짜 뉴스 수 : 17903 데이터 속 진짜 뉴스 수 : 21192

	title	text	subject	date
0	As U.S. budget fight looms, Republicans flip t	WASHINGTON (Reuters) - The head of a conservat	politicsNews	December 31, 2017
1	U.S. military to accept transgender recruits o	WASHINGTON (Reuters) - Transgender people will	politicsNews	December 29, 2017
2	Senior U.S. Republican senator: 'Let Mr. Muell	WASHINGTON (Reuters) - The special counsel inv	politicsNews	December 31, 2017
3	FBI Russia probe helped by Australian diplomat	WASHINGTON (Reuters) - Trump campaign adviser	politicsNews	December 30, 2017
4	Trump wants Postal Service to charge 'much mor	SEATTLE/WASHINGTON (Reuters) - President Donal	politicsNews	December 29, 2017
95	House panel chair introduces \$81 billion disas	WASHINGTON (Reuters) - The chairman of the U.S	politicsNews	December 19, 2017
96	Trump nominates Liberty University professor t	WASHINGTON (Reuters) - U.S. President Donald T	politicsNews	December 19, 2017
97	Trump on Twitter (Dec 18) - Congressional Race	The following statements were posted to the ve	politicsNews	December 18, 2017
98	Trump Cabinet officials to visit Puerto Rico t	WASHINGTON (Reuters) - Two members of Presiden	politicsNews	December 19, 2017
99	'Dreamer' immigration bill not on U.S. Senate	WASHINGTON (Reuters) - The U.S. Senate will no	politicsNews	December 18, 2017

출처:www.Kaggle.com

텍스트 데이터 전처리 과정

진짜 뉴스 데이터와 가짜 뉴스 데이터를 찾는다

찾은 데이터를 단어 단위로 토큰화한다

토큰화된 데이터에 존재하는 불용어들을 없앤다

불용어를 없앤 데이터의 어간을 추출하고 단어 원형을 복원한다

단어 원형이 복원된 토큰화된 데이터를 벡터화 시킨 후 학습한다

텍스트 데이터 전처리

텍스트 전처리 과정	개념	관련 개념/과정의 목 적	관련 개념 설명/과정의 효과
1. 텍스트 토 큰 화	장문의 데이터를 작은 단위로 나누는 작업	토큰/불용어 제거를 위해서	나뉜 문자열의 단위 (정하기에 따라 문장이나 단어, 또는 문자일 수 있음)
2. 불 용 어 제거	불용어인 단어들을 찾아 제거 하는 작업	불용어/문장에 있는 단어들 중 의미 없는 단어 들을 없애고, 의미 있는 단어들만 남기기 위해서	문장 내에서 많이 사용되지만, 문장의 전체 맥락과 상관없는 단어들 예)조사, 접속사, 접미사, 대명사 등
3. 어간 추출	단어에서 어미를 제거하고 어간을 추출하는 작업	같은 의미를 갖는 단어의 여러 가지 활용이 있을 경우 다른 단어로 인식되는 문제점 해결	문법적으로 변형된 단어들의 원형을 찾아 분석할 때 정확도를 높임
4. 표제어 추출	단어를 단어의 기본형으로 복원하는 작업	문법적 또는 의미적으로 변한 단어의 원형을 찾 기 위해서	품사와 같은 문법적 요소뿐만 아니라 문장 내에서 변한 단어들의 원형을 찾음, 어간 추출보다 더 정확 함

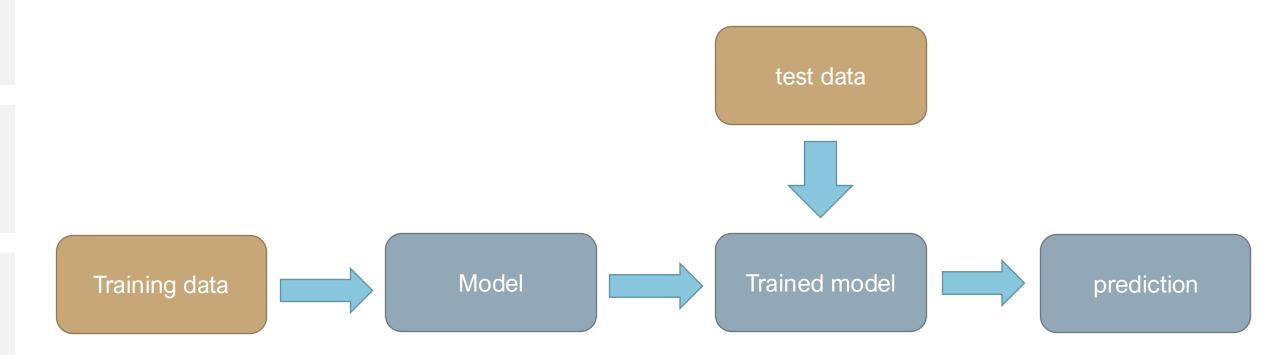
텍스트 데이터 전처리 결과

	original_news	preprocessed_news
0	WASHINGTON (Reuters) - The head of a conservat	[washington, reuters, head, conservative, repu
1	WASHINGTON (Reuters) - Transgender people will	[washington, reuters, transgender, people, all
2	WASHINGTON (Reuters) - The special counsel inv	[washington, reuters, special, counsel, invest
3	WASHINGTON (Reuters) - Trump campaign adviser	[washington, reuters, trump, campaign, adviser
4	SEATTLE/WASHINGTON (Reuters) - President Donal	[seattle, washington, reuters, president, dona
95	WASHINGTON (Reuters) - The chairman of the U.S	[washington, reuters, chairman, u, house, repr
96	WASHINGTON (Reuters) - U.S. President Donald T	[washington, reuters, u, president, donald, tr
97	The following statements were posted to the ve	[following, statement, posted, verified, twitt
98	WASHINGTON (Reuters) - Two members of Presiden	[washington, reuters, two, member, president,
99	WASHINGTON (Reuters) - The U.S. Senate will no	[washington, reuters, u, senate, consider, imm

이후, TF-IDF / Word2Vec과 같은 벡터화 과정을 거친 후, 모델 훈련 예정

Part 3 향후계획

Model Building & Training Framework

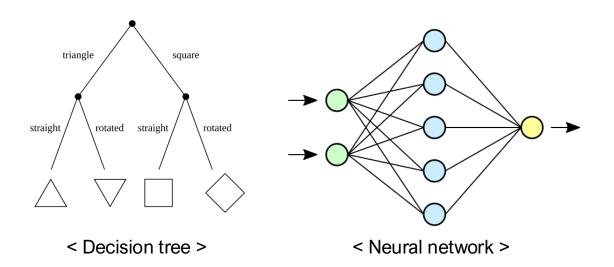


Model building & training

Model building

- 머신러닝 분류모델 구축

Ex) decision tree, neural network, support vector machine 등



Model training

- 1) 훈련, 테스트 데이터 3:1 비율로 나눔 (+ 훈련데이터의 일부를 검증데이터로 지정 → 모델을 일반화하는데 도움 줌)
- 2) 훈련데이터만으로 모델을 학습
- 3) 테스트 데이터로 모델의 성능을 수치화



Evaluation

		실제 정답	
		True	False
분류	True	True Positive	False Positive
결과	False	False Negative	True Negative

• True : 실제 정답과 분류 결과가 같음

• False : 실제 정답과 분류 결과가 다름

Positive : 분류 결과가 yes 또는 1

• Negative : 분류 결과가 no 또는 0

- 분류의 평가지표

$$(Recall) = \frac{TP}{TP + FN} \qquad (Accuracy) = \frac{TP + TN}{TP + FN + FP + TN}$$

$$(Precision) = \frac{TP}{TP + FP} \qquad (F1-score) = 2 \times \frac{1}{\frac{1}{Precision} + \frac{1}{Recall}} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

