# 13. 텍스트 마이닝(Text Mining)

# 1강. 웹 크롤링으로 기초 데이터 수집하기

### 학습내용

- 웹 크롤링으로 기초 데이터 수집하기
- 나무위키 최근 변경 페이지 키워드 분석

### 학습목표

- 웹 크롤링으로 기초 데이터를 수집하고 분석할 수 있다.
- 나무위키에서 최근 변경이 일어난 페이지들의 키워드를 이용하여 분석할 수 있다.

## 1. 웹 크롤링으로 기초 데이터 수집하기

- 1) 텍스트 마이닝
- 비정형 데이터, 텍스트 데이터로부터 유의미한 정보를 추출하는 데이터 분석
- 비정형 데이터는 이미지 데이터나 음성 데이터와 같은 정해진 형태가 없고 연산이 불가능한 데이터를 의미

#### 2) 텍스트 마이닝의 예

- 구글, 네이버 번역기, 챗봇
- 텍스트 임베딩(자연어 처리)
- 검색, 추천시스템
- 자연어 처리는 텍스트가 기계적으로 어떤 의미를 가지고 있는지를 추출하는 것

#### 3) 웹 크롤링

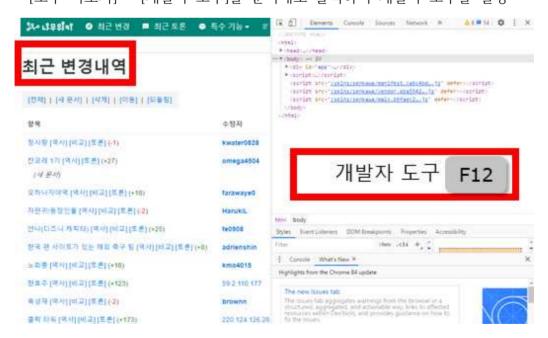
- 웹 스크래핑(Web Scraping)이라고도 하며 인터넷에 있는 웹 페이지를 방문해서 페이지의 자료를 자동으로 수집하는 작업
- 텍스트 데이터를 웹 크롤링으로 수집한 다음 데이터 내에서 등장한 키워드의 출현 빈도를 분석할 수 있음

### 가. 웹 크롤링으로 기초 데이터 수집하기

- 크롤링을 위한 첫 번째 단계는 인터넷 익스플로러, 크롬 등의 웹 브라우저를 실행하여 크롤링의 대상이 될 페이지 구조를 살펴보는 것



- [도구 더보기] - [개발자 도구]를 순서대로 클릭하여 개발자 도구를 실행



#### 나. 웹 크롤링 라이브러리

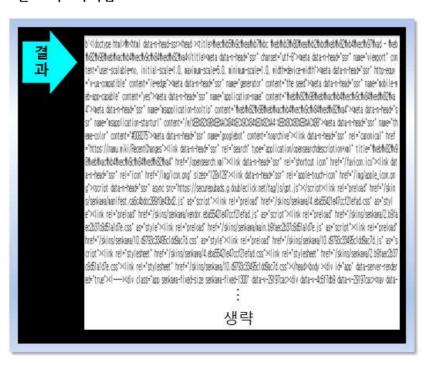
- BeautifulSoup과 requests 라는 라이브러리로 웹 크롤러를 만들 수 있음
- requests는 특정 URL로부터 HTML 문서를 가져오는 작업 수행
- BeautifulSoup 모듈은 HTML 문서에서 데이터를 추출하는 작업 수행

### 4) 페이지의 URL 정보 추출하기

- requests . get( ) 함수로 URL의 HTML 문서 가져옴

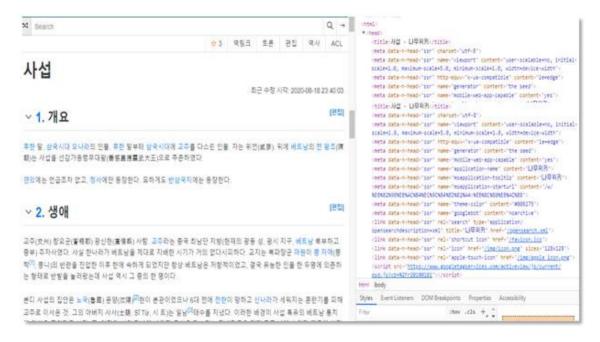


#### - html 코드를 모두 가져옴



- 파싱Parsing은 어떤 페이지(문서, html)에서 내가 원하는 데이터를 특정 패턴이나 순서로 추출하여 정보로 가공하는 것

- HTML 구조에 기반하여 table → tbody → tr → td → a 태그 순으로 HTML 계층 구조를 좁혀나가는 과정
- 목표 태그에 도달했을 때 get(herf)함수로 태그의 속성 정보를 추출
- get()함수는 해당 태그가 가지고 있는 특정한 속성을 추출
- 5) 텍스트 정보 수집하기
- 추출한 웹 페이지들의 URL을 방문하여 HTML 구조를 개발자 도구로 살펴 본 것



- get( ) 함수 대신 text( )함수를 사용하여 태그의 텍스트 정보만 추출

```
req = requests.get(page_urls[0])
html = req.content
soup = BeautifulSoup(html, 'lxml')
contents_table = soup.find(name="article")
title = contents_table.find_all('h1')[0]
category = contents_table.find_all('ul')[0]
content_paragraphs =
contents_table.find_all(name="div",
attrs={"class":"wiki-paragraph"})
content_corpus_list = []

print(title.text)
print(category.text)
print(content_corpus)
```

결 과 MA

한나간의 관료동오의 인물베트남의 군주추존된 왕137년 출생226년 사망

사업 가문의 역대 군주 베트남 중 지해 ← 초대 사업 → 2대 사회 陳朝 - 胡明 편 왕조 - 대우 황제 [ 聖치기 · 접기 ]陳朝廷 왕조세대 대중제2대 성종제3대인종제4대영종제5대영종제5대영종제7대유종제8대전폐제제9대영종홍큐제10대영종홍큐爾편 왕조胡明호 왕조胡明호 왕조대대우 폐제제12대순종제13대소재제1대국조제2대호환창映陳 후 편 왕조제14대간영제제15대중광제제16대진고추존 영무왕 · 목조 · 영조 · 원조 · 대조 · 충무왕 · 홍도왕 · 황대박 비정봉 진도에 · 진익작동라보기陳明한 왕조제1대단종제2대 성종제3대인종제4대영종제5대명종제6대현종제 7대유종제6대전폐제제9대여종蔣宗제10대영종寧宗陳明 전 왕조태明호 왕조대대1대학폐제제12대순종제13대소제제1대국조제2대호환창明陳明 후 전 왕조제14대간정제제15대중광제제16대건교추존 영무왕 · 목조 · 영조 · 원조 · 태조 · 충무왕 · 홍도왕 · 황대박 비정통 진도에 · 진익작 사호 선감가응령무대왕(善慈慕應建武大王) 왕호 사왕(土王, 시보엉)57 Veons) 성제 사(土) 휘 섭(葵) 베트남식 이름 시 니앱(57 Min lab) 생물 년도 137년 · 228년 제위 년도 187년 · 228년 1, 개요2, 성대3, 사후 4, 평가5, 미디어 막스5.1, 삼국제 시간조5.2, 삼국전투기5.3, 토 달 위: 삼국6, 물러보기후한 말, 삼국시대 오나라의 인물, 후한 말부터 삼국시대에 교주를 다스린 인물, 자는 위언(威彦), 위에 베트남의 전 왕 조(陳朝)는 사업을 선감기응경무대왕(書蔣蔣德建武大王)으로 추존하였다.연의에는 언급조차 없고, 정사에만 등장한다. 묘하게도 반삼국지에는 등 장한다.교주(兗州) 창오군(藩蔣화) 광신현(唐信朝) 사람, 교주라는 중국 최남단 지방(현재의 광웅 성, 광시 지구, 베트남 북부하고 중부) 주자사 였다. 사실 한나라가 베트남을 제대로 지배한 시기가 거의 없다시피하다. 교지는 복화장군 마원이 중 제대(종작[1], 종나)의 반안을 진압한 이후 하에 속하게 되었지만 한산 베트남은 제한적이었고, 결국 유남한 인물 한 두명에 의존하는 형태로 반발을 눌러왔는데 사선 역시 그 중의 한 명이

- 2. 나무위키 최근 변경 페이지 키워드 분석하기
- 1) 나무위키의 최근 데이터 크롤링하기
- 크롤링한 데이터를 데이터 프레임으로 만들기 위해 준비

```
columns = ['title', 'category', 'content_text']

df = pd.DataFrame(columns=columns)

# 각 페이지별 '제목', '카테고리', '본문' 정보를 데이터 프레임
으로 만듭니다.
for page_url in page_urls:

# 사이트의 html 구조에 기반하여 크롤링을 수행합니다.
  req = requests.get(page_url)
  html = req.content
  soup = BeautifulSoup(html, 'lxml')
  ::

df.head(5)
```

- 각 페이지별 제목, 카테고리, 본문 정보를 데이터 프레임으로 생성



### 2) 키워드 정보 추출

- 수집한 텍스트 데이터에서 키워드 정보를 추출하기 위해 텍스트 전처리 작업이 필요
- 텍스트 전처리는 특수문자나 외국어를 제거하는 등의 과정을 포함
- 키워드 정보 추출은 좁은 의미에서
- 명사 혹은 형태소 단위의 문자열을 추출하는 것

#### 가. 키워드 정보 추출에서 언어와 상황마다 다른 경우

- 스팸메일을 분류하는 텍스트 마이닝의 경우 특수문자나 외국어가 분석의 중요 한 힌트가 될 수 있기 때문에 이를 제거하지 않는 편임
- 키워드 분석처럼 '단어'를 추출하는 것이 목적이라면 특정 언어의 글자만을 추출하기도 함

### 나. 정규표현식

- 정규표현식이란 특정 규칙을 가진 문자열의 집합을 표현하는 형식
- 파이썬에서는 're'라는 모듈을 통해 정규표현식을 사용
- '[^ ¬ ] 가 힣]+' 이라는 코드로 한글에 대한 정규표현식을 정의하면 대 상이 되는 텍스트 데이터에서 한글만 추출할 수 있음

### 3) 텍스트 데이터 전처리하기

```
# 텍스트 정제 함수 : 한글 이외의 문자는 전부 제거합니다.

def text_cleaning(text):

hangul = re.compile('[^ ¬- | 가-헿]+')

# 한글의 정규표현식을 나타냅니다.

result = hangul.sub('', text)

return result

print(text_cleaning(df['content_text'][0]))
```



4) 모든 데이터에 전처리하기

```
# 각 피처마다 데이터 전처리를 적용합니다.

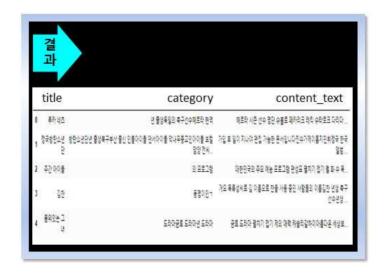
df['title'] = df['title'].apply(lambda x: text_cleaning(x))

df['category'] = df['category'].apply(lambda x:
text_cleaning(x))

df['content_text'] = df['content_text'].apply(lambda x:
text_cleaning(x))

df.head(5)
```

- 모든 데이터에 전처리를 적용하기 위해서는 apply()함수 사용



# 평가하기

- 1. 인터넷에 있는 웹 페이지를 방문해서 페이지의 자료를 자동으로 수집하는 작업은?
  - ① 핫키
  - ② 리스트 URL
  - ③ 웹 크롤링
  - ④ 웹 판다스
- 정답 : ③번

해설: 인터넷에 있는 웹 페이지를 방문해서 페이지의 자료를 자동으로 수집하는 작업을 의미합니다.

- 2. 웹 크롤링 라이브러리 중 특정 URL로부터 HTML 문서를 가져오는 작업을 수행하는 것은?
  - ① BeautifulSoup
  - ② Requests
  - 3 Numpy
  - Pandas
- 정답 : ②번

해설 : Requests는 특정 URL로부터 HTML 문서를 가져오는 작업을 수행합니다.

### 학습정리

- 1. 웹 크롤링으로 기초 데이터 수집하기
  - 텍스트 마이닝
  - 웹 크롤링의 개념, 라이브러리
- 2. 나무위키 최근 변경 페이지 키워드 분석하기
  - 웹 데이터 가져오기
  - 키워드 정보 추출하기