

분석용 데이터 획류젤략

웹 리뷰 데이터 확보 기법

학습 목표

+ + +

학습 목표

- 결과를 저장할 폴더와 파일명을 자동으로 생성하도록 설정할수 있다.
- 리뷰 데이터에서 원하는 항목을 선택하여 추출할 수 있다.
- 다양한 사이트별 리뷰 데이터 수집 방법을 활용할 수 있다.

학습 내용

- 다양한 형태의 리뷰 데이터 수집하기 이론
- 다양한 형태의 리뷰 데이터 수집하기 실습

인터넷 뉴스 리뷰 정보 수집하기 - 이론

1) 리뷰 수집용 인터넷 뉴스 예제



인터넷뉴스예제



다양한 형태의 리뷰 데이터 수집하기 - 이론

인터넷 뉴스 리뷰 정보 수집하기 - 이론

- 2) 총 리뷰 건수 정보 수집
 - (1) 작업 순서
 - 구성 및 변경 관리
 - ① 사용자에게 수집할 리뷰 건수 물어보기
 - ② 결과를 저장할 폴더명 물어보기
 - 의 에이터를입력받은 후검색을시작하여
 - ▶ 전체 검색 결과 건수
 - ᆗ 실제최종출력건수
 - ➡ 총검색페이지수정보를보여주고수집시작

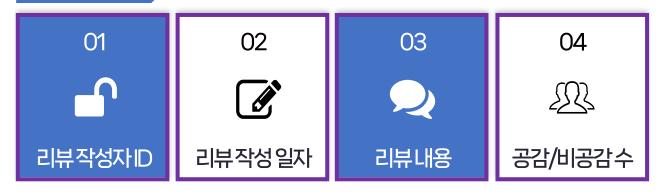
주의사항

사용자가요청한건수와실제리뷰건수와비교하여작업조율이필요

인터넷 뉴스 리뷰 정보 수집하기 - 이론

3) 수집할 상세 항목 정보





- 4) 추가 리뷰 정보 보기
 - 네이버 뉴스 리뷰
 - ☑ 맨처음페이지는5개리뷰만
 - '댓글더보기>'클릭하면추가리뷰정보
 - **I** 1페이지에 20개의리뷰
 - 이후리뷰를더보려면이래'더보기'버튼클릭

다양한 형태의 리뷰 데이터 수집하기 - 이론

인터넷 뉴스 리뷰 정보 수집하기 - 이론

- 5) 리뷰 정보 저장하기
- 아래와같이 폴더와 파일을 자동으로 생성함
 - + 자동생성폴더이름형식
 - 자동생성파일이름형식

년도4자리-월-일-시간-분-초-뉴스기사댓글.csv 년도4자리-월-일-시간-분-초-뉴스기사댓글.txt 년도4자리-월-일-시간-분-초-뉴스기사댓글.xls

주의사항

파일자동생성시파일이름 중복으로 덮어쓰는 문제가발생

다양한 형태의 리뷰 데이터 수집하기 - 실습

인터넷 뉴스 리뷰 정보 수집하기 - 실습

■ 검색어 지정

```
1 #Step 1. 필요한 모듈과 라이브러리를 로딩합니다.
2 from bs4 import BeautifulSoup
3 from selenium import webdriver
4 from selenium.webdriver.common.by import By
  from selenium.webdriver.common.keys import Keys
6 from selenium.webdriver.chrome.service import Service
7 import time
8 import math
9 import numpy
10 import pandas as pd
   import os
12
13
  #Step 2. 사용자에게 검색어 키워드를 입력 받고 저장할 폴더와 파일명을 설정합니다.
14
  print("=" *80)
15
  print("뉴스 기사의 댓글 정보 수집하기")
16
   print("=" *80)
                                                                           I
   query txt = '뉴스기사댓글'
  query_url = 'https://news.naver.com/main/read.nhn?mode=LSD&mid=shm&sid1=102&oid=056&aid=0010661268
21
      cnt = int(input('1.크롤링 할 뉴스 리뷰 건수는 명건입니까?(기본값: 20건): '))
23
  except ValueError :
24
      cnt = 20
25
26
  page cnt = math.ceil(cnt / 20)
27
28
  f_dir = input("2.결과를 파일을 저장할 폴더명만 쓰세요(기본값:c:\\py_temp\\):")
29 if f_dir=''
30
      f_dir='c:\\py_temp\\'
```

주의사항

리뷰건수를입력안하고실행할경우 ValueError 발생

인터넷 뉴스 리뷰 정보 수집하기 - 실습

■ 저장할 폴더 지정

```
23 except ValueError:
24
       cnt = 20
25
   page cnt = math.ceil(cnt / 20)
26
28 |f_dir = input("2.결과를 파일을 저장할 폴더명만 쓰세요(기본값:c:\\py_temp\\):")
   if f dir="
29
       f_dir='c:\\py_temp\\
30
   # 저장될 파일위치와 이름을 지정합니다
32
   now = time.localtime()
33
   s = \frac{804d-802d-802d-802d-802d-802d}{0.0000} (now.tm_year, now.tm_mon, now.tm_mday, \
35
                                            now.tm_hour, now.tm_min, now.tm_sec)
36
   os.makedirs(f_dir+s+'-'+query_txt)
37
38 os.chdir(f_dir+s+'-'+query_txt)
39
40 ff_name=f_dir+s+'-'+query_txt+'\\'+s+'-'+query_txt+'.txt'
41 fc_name=f_dir+s+'-'+query_txt+'\\'+s+'-'+query_txt+'.csv'
   fx_name=f_dir+s+'-'+query_txt+'\\'+s+'-'+query_txt+'.xls'
42
43
   #Step 3. 크롬 드라이버를 사용해서 웹 브라우저를 실행합니다.
44
   s time = time.time( )
45
46
   s = Service("c:/py_temp/chromedriver.exe")
47
48 driver = webdriver.Chrome(service=s)
49
50 driver.get(query_url)
   driver.maximize_window()
52 time.sleep(5)
5.2
```

주의사항

폴더/파일이름에 <mark>현재날짜, 시간등을 넣으면 덮어쓰지 않고</mark> 저장가능

정보

time.localtime()함수:현재시간을년,월,일,시등의형태로보여주는함수

주의사항

time.time()함수:현재시간을실수형태로반환하여보여주는함수

인터넷 뉴스 리뷰 정보 수집하기 - 실습

■ 콤마(,) 기호 삭제

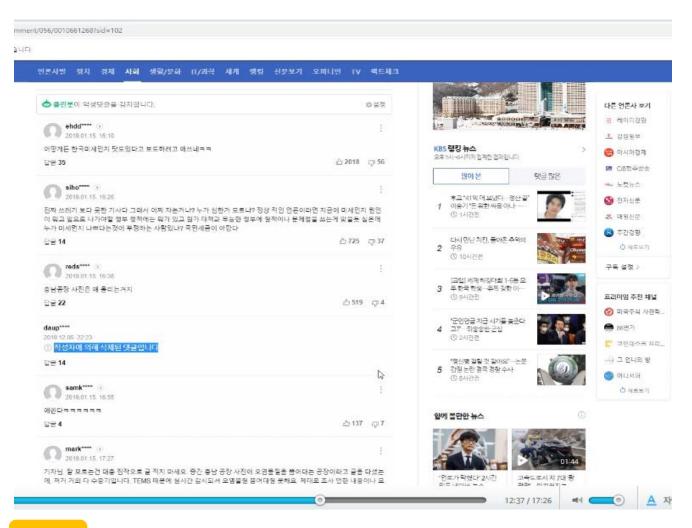
```
48 driver = webdriver.Chrome(service=s)
49
50 driver.get(query_url)
51 driver.maximize_window()
   time.sleep(5)
    # 현재 총 리뷰 건수를 확인하여 사용자의 요청건수와 비교 후 동기화합니다
    html = driver.page_source
   soup = BeautifulSoup(html, 'html.parser')
58 result= soup.find('div','media_end_head_info_variety_cmtcount _COMMENT_HIDE').get_text()
59
60 print("=" *80)
51 search_cnt = int( result.replace(",
63 if cnt > search_cnt :
         cnt = search_cnt
64
65
66 print("전체 검색 결과 건수 :",search_cnt,"건")
67 print("실제 최종 출력 건수 :",cnt)
68 print("총 페이지 수:" , page_cnt)
 1 # 사용자가 요청한 건수가 많을 경우 리뷰 더보기 버튼을 클릭합니다
2 # 최초 10건 수집후 댓글 더보기 버튼 클릭
3 # 아레 버튼을 눌러 첫 화면에 총 20건의 댓글이 나오게 만듦
4 driver.find_element(By.XPATH,'//*[@id="cbox_module"]/div/div[9]/a/span[1]').click()
   time.sleep(3)
7 #Step 6, 20건 출력되어 있는 현재 페이지 리뷰와 점수 등 내용 수집
8 no2=[] # 리뷰 번호
9 writer_id2=[] # 리뷰 작성자 ID
10 review2=[] # 리뷰 작성 일자
11 write_date2=[] # 리뷰 작성 일자
12 nogam 0=[] # 곳간 회수
 6
10 review2=[]
11 write_date2=[]
12 gogam 0=[]
```

정보

cnt:사용자가입력한값, search_cnt:실제 검색된 건수

다양한 형태의 리뷰 데이터 수집하기 - 실습

인터넷 뉴스 리뷰 정보 수집하기 - 실습



정보

작성자에의해삭제된댓글은 공감/비공감표시가없음

인터넷 뉴스 리뷰 정보 수집하기 - 실습

```
3 # 아래 버튼을 눌러 첫 화면에 총 20건의 댓글이 나오게 만듦
4 driver.find_element(By.XPATH,'//*[@id="cbox_module"]/div/div[9]/a/span[1]').click()
 5 time.sleep(3)
 6
   #Step 6. 20건 출력되어 있는 헌제 페이지 리뷰와 점수 등 내용 수집 no2=[] # 리뷰 번호 writer_id2=[] # 리뷰 작성자 ID review2=[] # 리뷰 내용 write_date2=[] # 리뷰 작성 일자 gogam_0=[] # 공감 횟수 gogam_1=[] # 비공감 횟수
8 no2= []
 9 writer_id2=[]
10 review2=[]
13 gogam_1=[]
14 count = 0
15
16 for a in range(1,page_cnt+1) :
17
18
        if a == page_cnt :
                break
19
        else :
20
21
             driver.find_element(By.XPATH)'//*[@id="cbox_module"]/div/div[9]/a').click()
22
             time.sleep(3)
23
             print("%s페이지 이동 완료=
             time.sleep(random.randrange(1,3)) # 3-8 초 사이에 랜덤으로 시간 선택
24
25
26
   print('이제 리뷰 정보를 수집합니다. 잠시만 기다려 주세요~~~~~')
28 #txt 파일에 저장하기 위해 파일 open하기
29 f = open(ff_name, 'a',encoding='UTF-8')
30
31
   html = driver.page_source
   soup = BeautifulSoup(html, 'html.parser')
32
33
   slist = soup.find('ul', 'u_cbox_list').find_all('li')
34
36 for li in slist:
```