



INTRODUCTION

프렌차이즈순위에서도 커피창업은 줄곧 상위를 지키고 있습니다. 프렌차이즈 창업의 커피 (디저트)의 창업율은 40%에 달한다고 합니다. 통계청에 의하면 우리나라 사람 1인당 연간 커피 소비량은 484잔이라고 합니다. 그럼 매일 1인당 한 잔 이상씩 마시고 있는 것과 같습니다.

이처럼 빠른 성장세를 보이고 있는 커피시장에서 고객을 사로잡을 수 있는 확률이 높은 신제품을 제안하기 위해 이 프로젝트를 시작하게 되었습니다.

TABLE OF CONTENTS

PART 1. 네이버 블로그

PART 2. 인스타그램













PART 1. 네이버블로그

- I .데이터 수집
- Ⅱ.텍스트 정제
 - 1. word2vec 알고리즘
 - 2. 단어 통일
 - 3. 게시글 정제
 - 4. 기타 정제
- Ⅲ.텍스트 정보화
 - 1. 게시글 점수화
 - 2. 점수 테이블 만들기
 - 3. 점수 부여

Ⅳ. 텍스트 학습

- 1. 데이터셋 만들기
- 2. 데이터 병합
- 3. 점수 라벨링
- 4. 머신러닝

Ⅳ. 결과

네이버 블로그 크롤링이 안된다...?



"selenium"을 통한 게시글 크롤링

─ 웹 브라우저를 컨트롤하여 UI(User Interface)를 지정하는 도구

```
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
driver = webdriver.Chrome("c:/Windows/chromedriver")
for q in starbucks['click']:
    except_starbucks_naver = []
    driver.get("http://m.naver.com")
    driver.find_element_by_css_selector('#query').send_keys(q)
    driver.find_element_by_css_selector('#query').send_keys('\n')
......

for z in range(1,11):
    while i < 16:
        num = (page-1)*15
        try:
        time.sleep(random.randrange(1,2))
        y = driver.find_element_by_css_selector('#blog_' +str(num+i)+' > a > div.total......
```

스타벅스 블로그 : 12,424개

+ 이디야 블로그 : 5,607개

합계 : 18,031개

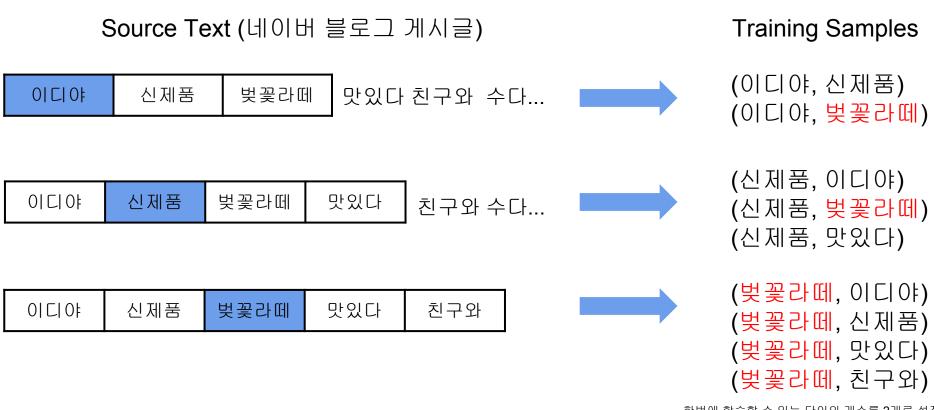
18,031개.. 블로그 어떻게 정제작업을 할까...?

II - 1. Word2Vec 알고리즘



"word2vec"을 활용한 단어간 연관성분석

단어들을 벡터화시켜 단어간의 거리와 말뭉치의 빈도수를 계산해주는 알고리즘



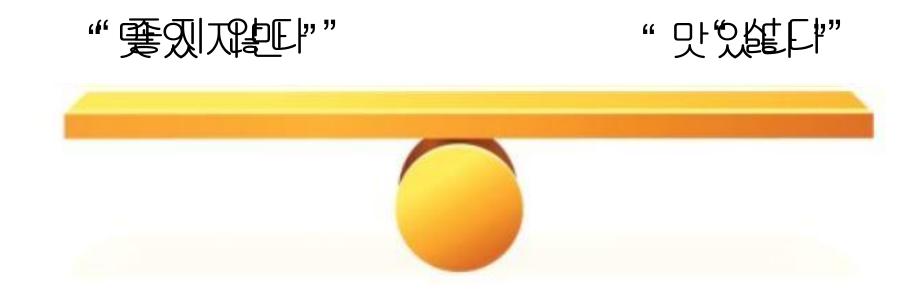
한번에 학습할 수 있는 단어의 개수를 2개로 설정

Ⅲ - 2. 단어 통일



"말이 아 다르고 어 다르다"

고객의 미세한 감정을 잘 활용해야 한다.



Ⅲ - 2. 단어 통일



선행부정어

안

후행부정어



않, **없**, 아니, 모르

반대 내용을 이끄는 접속부사



하나, 지만, 한데, 하지만, 그러나

기타 처리



입맛에 딱, ~와 똑같다

Ⅲ - 2. 단어 통일



안맛있다

맛있지않다

" 사전등록 "



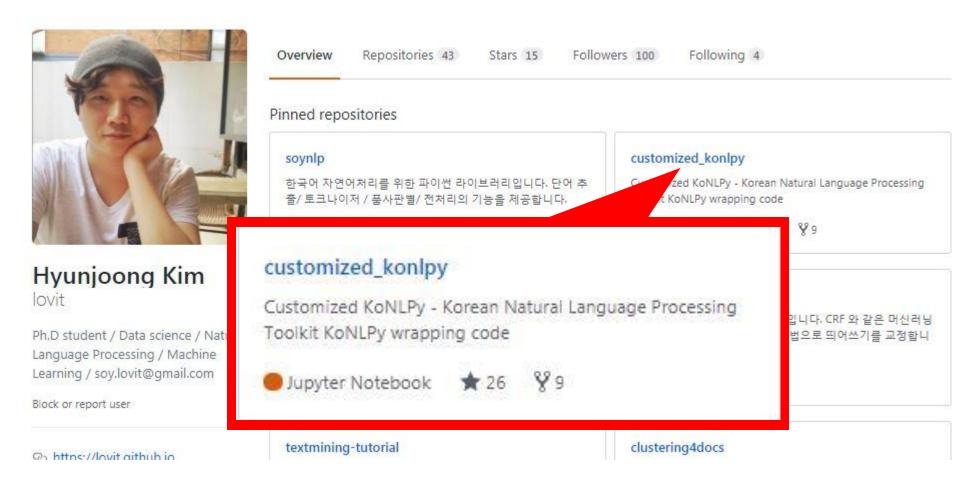
(벚꽃라떼,안맛있다) (벚꽃라떼, 맛있지않다) (벚꽃라떼, 맛있지만) (벚꽃라떼, 맛있다)

맛있지만

맛있다

Ⅲ - 3. ckonlpy





https://github.com/lovit

Ⅲ - 3. 게시글 정제 작업





나는 이 음료에 시시한 추억을 가지

2009년 겨울, 내 생에 첫 텀블러 쿠폰으로 주문한 음료가 년

~/r 프라푸치노인 것이다.

그 당시 텀블러 쿠폰은 그야말로 무리 ㅋ 쿠폰이었다.

음료 종류와 사이즈 불문, 엑스트라 횟수 제한도 없어 뭐든지 가능했다.

물론 스타벅스의 엑스트라 서비스를 잘 아는 사람도 없을 때였다.

어쨌든 몹시 다정한 파트너가 조심스레 포장한 프라푸치노를 들고 총총걸음으로 기숙사에 돌아가던 길은 설레고 즐거웠다.

추고 매서운 바란이 분어도 마냥 시난다



```
re adjective('굴','좋다','좋지않다','좋지만',docs)
re adjective('색다','색다르다','색다르지않다','색다르지만',docs)
re adjective('다르','다르다','다르지않다','다르지만',docs)
re adjective('똑같','똑같다','똑같지않다','똑같지만',docs)
re adjective('귀여','귀엽다','귀엽지않다','귀엽지만',docs)
re adjective('아까','아깝다','아깝지않다','아깝지만',docs)
re adjective('괜찮','괜찮다','괜찮지않다','괜찮지만',docs)
re adjective('따똣','따뜻하다','따뜻하지않다','따뜻하지만',docs)
re adjective('나쁘','나쁘다','나쁘지않다','나쁘지만',docs)
re adjective('아쉬','아쉽다','아쉽지않다','아쉽지만',docs)
re adjective('아쉽','아쉽다','아쉽지않다','아쉽지만',docs)
re noun('의망','의망하다','의망하지않다','의망하지만',docs)
re noun('달달','달달하다','달달하지않다','달달하지만',docs)
re noun('달','달달하다','달달하지않다','달달하지만',docs)
re noun('뿔뿔','뿔뿔하다','뿔뿔하지않다','뿔뿔하지만',docs)
re noun('실패','실패하다','실패하지않다','실패하지만',docs)
re noun('추천','추천하다','추천하지않다','추천하지만',docs)
re noun('취전','취전하다','취전하지않다','취전하지만',docs)
re noun('실쿵','실쿵하다','실쿵하지않다','실쿵하지만',docs)
re noun('합격','합격하다','합격하지않다','합격하지만',docs)
re noun('풍미','풍미롭다','풍미롭지않다','풍미롭지만',docs)
no noun(' ) = ' ' ] = T(C)' ' ] = T(O)(C)' ' ] = T(T(D)' does)
```



```
'입맛에 딱이였음!!!!!! 아',
'입맛은 아니라며 복숭아',
'입맛 무엇..?..) 카페에서 쇼케이스를',
'입맛인지. 그래도 이제',
'입맛에는 아니예요. ㅋ',
'입맛돋고 좋드라구요^^ 솨고기',
'입맛에 너무 달다ㅠ',
                   カ시인.. 합심인. . 그래까 흑싙 ,
'입맛인 나에게 정말',
'입맛에는 괜찮았어여!! 그런데 '그 헤이즐럿 아메리카노와 똑같',
                   '엑스트라로 카라멜드리즐을 추가했다. 똑같',
'입맛에딸!!맞느ㅎㅎ 살짝'
                   '에이 이거 딸기우유랑 똑같',
                   '항아리모양 딸기우유랑 맛 똑같',
                   '오 블로그후기를에서 봤던거랑 똑같']
```

₩ - 1. 게시글 점수화



```
from ckonlpy.tag import Twitter
t = Twitter()
final beverage = []
postprocessor = Postprocessor(t, replace = dic replace)
file name kobill = kobill.fileids()#kobill 약에 있는 파일이름
docs ko = [kobill.open(i).read() for i in kobill.fileids()]
from gensim.models import word2vec
for i in range(len(docs ko)):
    print(file name kobill[i])
    if len(re.findall('\w{0,}',docs ko[i])) != len(docs ko[i])+1:
        docs ko[i] = docs ko[i]*20
        pos = lambda d: ['/'.join(p) for p in postprocessor.tag(d)]
        texts ko = pos(docs ko[i])
        for Z in range(len(texts ko)):
            texts ko[Z] = [texts ko[Z]]
       wv model ko = word2vec.Word2Vec(texts ko)
       wv model ko.init sims(replace=True)
       wv model ko.save('ko word2vec e.model')
        try:
            file name = file name kobill[i]
            a = score beverage()
            final beverage.append(a)
        except:
            continue
final = []
final.append(final beverage)
```

```
IPython console
Console 1/A 🔯
                                          ■ 8 Q
2018 및 꽃라떼 274. txt
('달달하다/Noun', -0.01860680803656578)
2018벚꽃라떼275.txt
2018벚꽃라떼 276. txt
('지다/Noun', 0.17207071185112)
('새롭다/Noun', 0.03664127737283707)
('저격하다/Noun', 0.02626614272594452)
('고급지다/Noun', 0.003623083233833313)
('아쉽다/Noun', 0.0029624737799167633)
('비싸다/Noun', -0.005824774503707886)
('달달하다/Noun', -0.01860680803656578)
('예쁘다/Noun', -0.024207476526498795)
('괜찮다/Noun', -0.047976039350032806)
('작다/Noun', -0.22834636270999908)
2018벚꽃라떼277.txt
2018벚꽃라떼 278.txt
('아쉽다/Noun', 0.0029624737799167633)
('예쁘다/Noun', -0.024207476526498795)
('작다/Noun', -0.22834636270999908)
2018벚꽃라떼279.txt
('예쁘다/Noun', -0.024207476526498795)
  . . . . . . . .
```

₩ - 2. 점수 테이블 생성



NEGATIVE



POSITIVE

-20점	-15점	-10점	-5점	+5점	+10점	+15점	+20점
좋지않다	싫지만	편하지않다	불편하지만	편하지만	편하다	좋지만	좋다
싫다	아깝다	불편하다	예쁘지만	괜찮지만	불편하지않다	싫지않다	굉장하다
맛없다	아쉽다	예쁘지않다	가깝지만	따뜻하지만	예쁘다	아깝지않다	추천하다
굉장하지않다	멀다	아깝지만	아름답지만	새롭지만	괜찮다	아쉽지않다	애정하다
나쁘다	맛있지않다	아쉽지만	뜨겁지만	적당하지만	따뜻하다	멀지않다	뿜뿜하다
추천하지않다	맛없지만	괜찮지않다	똑같지만	즐겁지만	새롭다	맛있다	저격하다
애정하지않다	부드럽지않다	따뜻하지않다		멋지지만	가깝다	부드럽다	심쿵하다
욕하다	시원하지않다	새롭지않다		많지만	맛있지만	시원하다	취저하다
뿜뿜하지않다	깔끔하지않다	가깝지않다		뜨겁지않다	맛없지않다	깔끔하다	추천하다
저격하지않다	고소하지않다	멀지만		귀엽지만	부드럽지만	고소하다	
심쿵하지않다	달콤하지않다	아름답지않다		재밌지만	아름답다	달콤하다	
취저하지않다	진하지않다	싱겁지만		설레지만	시원하지만	진하다	
추천하지않다	싱겁다	적당하지않다		만끽했지만	깔끔하지만	싱겁지않다	
	저렴하지않다	저렴하지만		반했지만	고소하지만	굉장하지만	
	비싸다	비싸지만			달콤하지만	저렴하다	
	비슷하다	비슷하지만			진하지만	비싸지않다	



def score beverage():

```
for i in wv model ko.most similar(pos(j),topn=2000):
    if (i[0][i[0].find('/')+1:] == 'Noun') and (i[0][:i[0].find('/')][-1:] == '<math>\mathcal{L}'):
        print(i)
        lst.append(i)
for i1 in wv model ko.most similar(pos(j),topn=2000):
    if (i1[0][i1[0].find('/')+1:] == 'Noun') and (i1[0][:i1[0].find('/')][-1:] == '지만'): {'카라멜리스트레토비안코82.txt': 1.9731426797807217},
        print(i1)
        lst.append(i1)
for i2 in wv model ko.most similar(pos(j),topn=2000):
    if (i2[0][i2[0].find('/')+1:] == 'Noun') and (i2[0][:i2[0].find('/')][-1:] in noun):
        print(i2)
        lst.append(i)
for z in range(len(lst)):
    if lst[z][0][:lst[z][0].find('/')] in positive 5:
        score = lst[z][1]
        score = score * 5
        score final += score
   if lst[z][0][:lst[z][0].find('/')] in negative 20:
       score = lst[z][1]
       score = score * (-20)
       score final += score
score dic[j1] = score final
return score dic
```

print score dic():

```
{'카라멜리스트레토비안코71.txt': 0.1971839740872383},
{'카라멜리스트레토비안코74.txt': -1.0110549442470074},
{'카라멜리스트레토비안코77.txt': -0.16882652416825294},
{'카라멜리스트레토비안코79.txt': -0.25920677930116653},
{'카라멜리스트레토비안코8.txt': -0.4964792914688587},
{'카라멜리스트레토비안코83.txt': 0.7857213355600834},
{'카라멜리스트레토비안코84.txt': 0.06514834240078926},
{'카라멜리스트레토비안코85.txt': 2.914459779858589},
{'카라멜리스트레토비안코87.txt': 0.9093738161027431},
{'카라멜리스트레토비안코90.txt': 1.02241612970829},
{'카라멜리스트레토비안코91.txt': 6.798028033226728},
{'카라멜리스트레토비안코92.txt': -0.5029981397092342},
{'카라멜리스트레토비안코94.txt': 1.2897718138992786},
{'카라멜리스트레토비안코95.txt': -1.5630182810127735},
{'카라멜리스트레토비안코96.txt': -0.4891696572303772},
{'카라멜리스트레토비안코97.txt': 1.6434356570243835},
{'벚꽃라떼1.txt': -2.7710622549057007},
{'벚꽃라떼10.txt': -0.19313422963023186}.
```



카라멜 소스 카라멜 드리즐 카라멜 카라멜 시럽 비니스트마일드 에스프레소 샷 샷 비니스트오리지널 과테말라 에스프레소샷

Ⅳ - 1. RECIPE 데이터셋 만들기



menu	샷	우유	 카라멜시 럽		카라멜드리즐	 1	카라멜소스	 벚꽃파우더	
카라멜리스트 레토비안코	1	0	 3	약 30	OO개에 가까{ column		0	 0	
			 		Column			 •••	
벚꽃라떼	0	1	 0		0		0	 1	
	•••	•••	 		•••			 •••	

menu	샷	우유		카라멜	 블라썸	
카라멜리스트레토비안코	1	0		1	 0	
	•••		119개의	column	 •••	
벚꽃라떼	0	1		0	 1	
	•••	•••			 •••	



RECIPE TABLE

menu	샷	우유	 점수
카라멜리스트레토비안코	1	0	
	•••		
벚꽃라떼	0	1	
	•••	•••	

SCORE_DIC()

{'카라멜리스트레토비안코1.txt': -0.5459930375218391},

•••

{'카라멜리스트레토비안코11.txt': 2.0612977258861065},

• • •

{'벚꽃라떼114.txt': -2.256549783051014}

•••

{'벚꽃라떼282.txt': 1.6226249374449253}



RESULT_TABLE

index	menu	샷	우유		점수
0	공주보늬밤라떼	1	1		
••••	•••	•••	•••		****
1700	벚꽃라떼	0	1		0
1701	벚꽃라떼	0	1		-2.5356519
	****	•••	•••	•••	••••
13121	카라멜리스트레토비안코	1	1		5.6823108
13122	카라멜리스트레토비안코	1	1		-1.2097426
	•••				
18030	화이트코코리스트레토	1	1		

Ⅳ - 3. 점수 라벨링



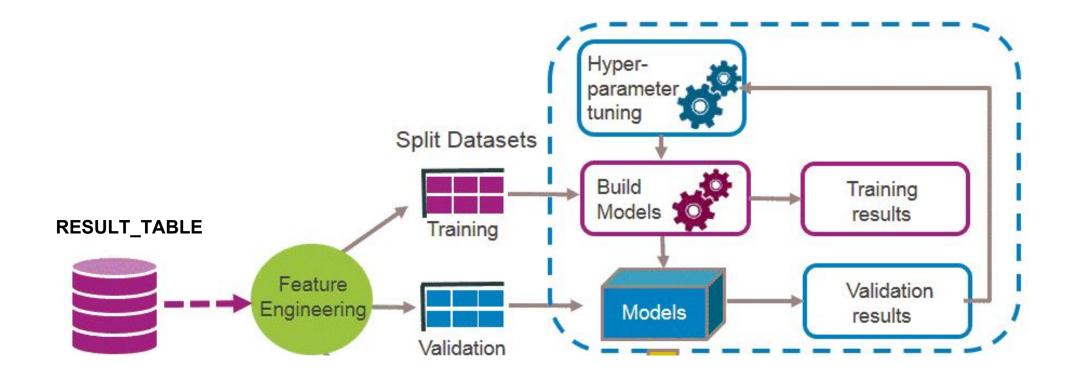
RESULT_TABLE

index		LABEL
0	•••	
1700		2
1701		1
13121		4
13122		1

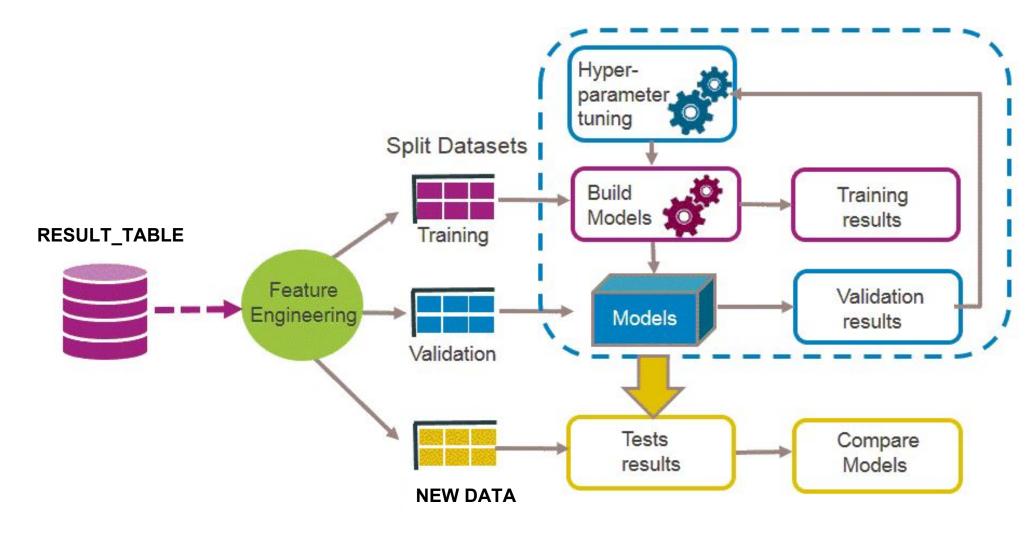
기준	LABEL
-10 이상 ~ - 5 미만	0
-5 이상 ~ 0 미만	1
0점	2
0점 초과 ~ 5점 이하	3
5점 초과 ~ 10점 이하	4

IV - 4. Machine Learning









IF, 딸기와 돌체시럽을 넣은 신제품이 나오면 어떤 점수가 나올까?



PART 2. 인스타그램

I. 데이터 수집

- 1. 문제점
- 2. 데이터 수집
- 3. 데이터 수집 코드
- 4. 데이터

Ⅱ. 데이터 정제

- 1. 인스타그램 약어 정제
- 2. 불필요한 단어 정제
- 3. ~그램 정제
- 4. 한글화
- 5. 기타 정제

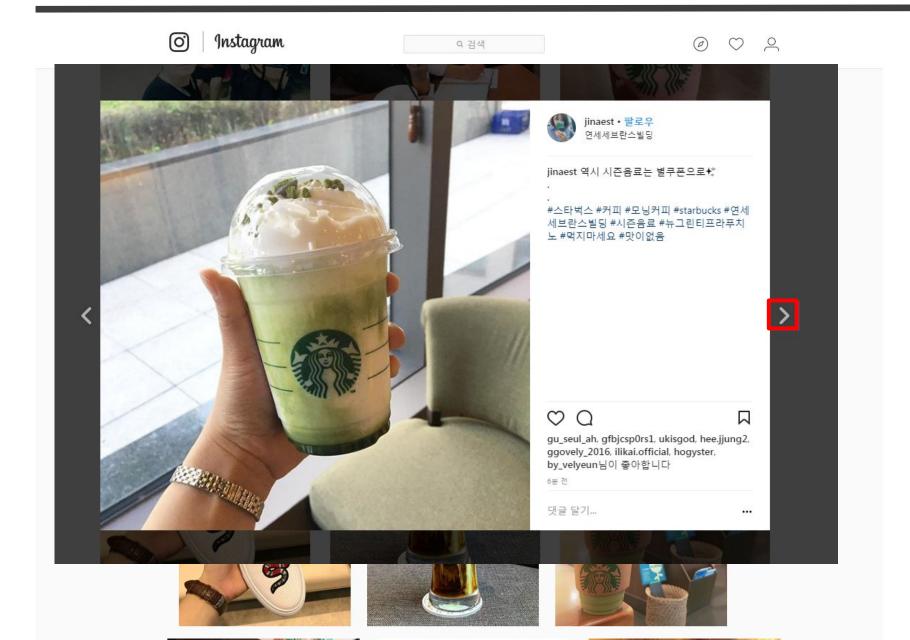
III. WORD CLOUD

- 1. 전체 데이터
- 2. menu 제외
- 3. 결과

₩. 이벤트 제안

I. 데이터 수집





I - 1. 문제점

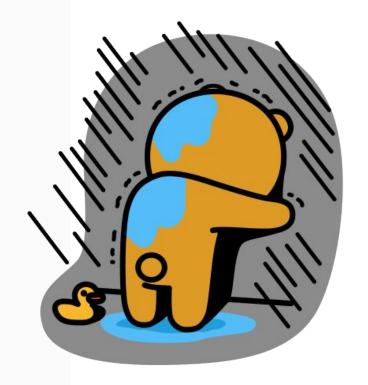


Instagram

죄송합니다. 페이지를 사용할 수 없습니다.

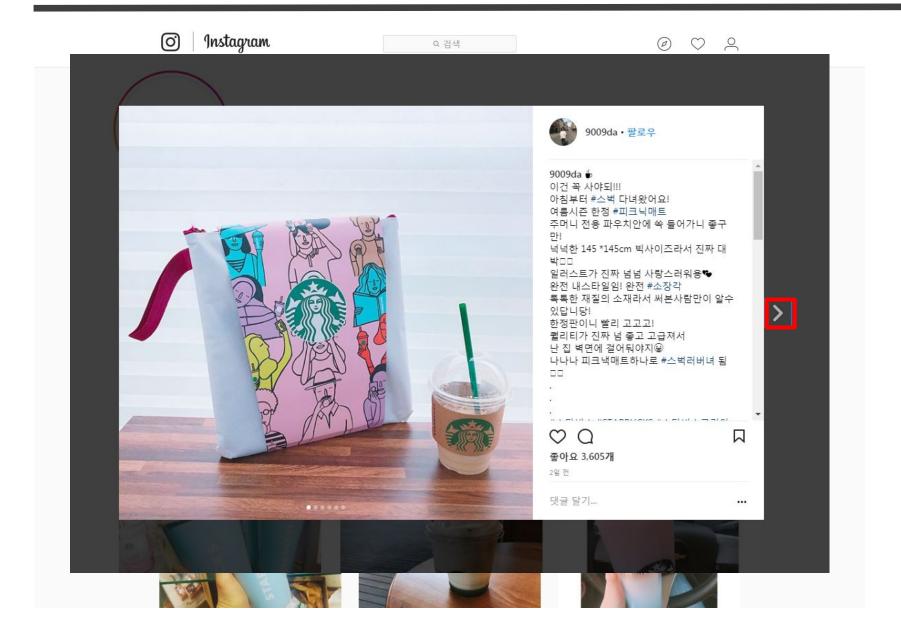
클릭하신 링크가 잘못되었거나 페이지가 삭제되었습니다. Instagram으로 돌아가기.

사라진 게시글 多



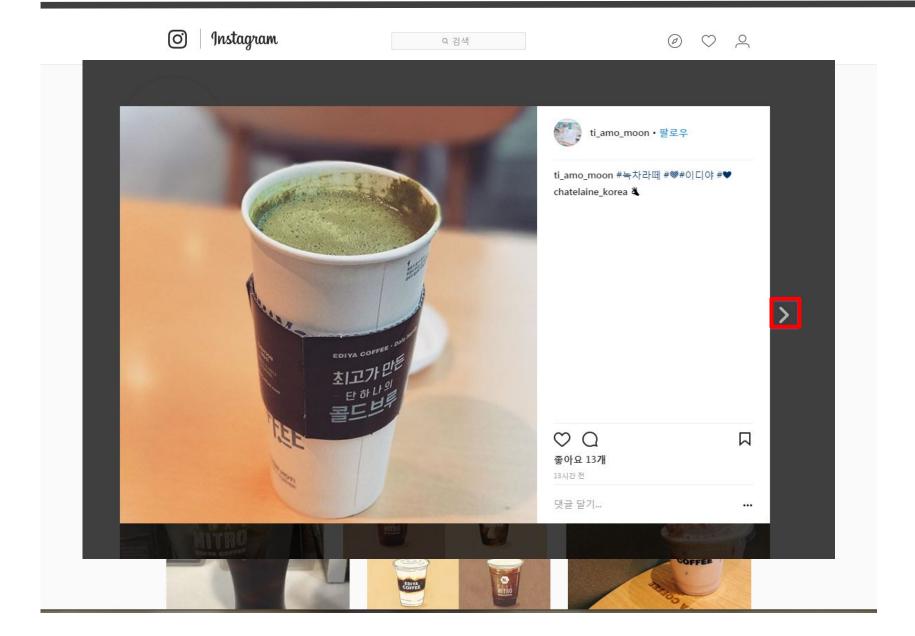
I - 2. 데이터 수집





I - 2. 데이터 수집





I - 3. 데이터 수집 코드



```
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
driver = webdriver.Chrome('C:\selenium\chromedriver.exe')
driver.get('http://www.instagram.com')
driver.find element by css selector('#email').send keys('****@*****.***')
driver.find element by css selector('#pass').send keys('******')
driver.find element by_css_selector('#loginbutton').click()
for i in ['이디야']:
       driver.find_element_by_css_selector('....').send_keys(i)
       time.sleep(random.randrange(1,3))
       driver.find element by css selector('....').send keys('\n\n')
       text1 = []
       if count != 0:
               driver.find element by css selector('....').click()
               for j in range(1, int(count)):
                      try
                            text1.append(....').text)
                      except Exception as error:
                             time.sleep(random.randrange(1,3))
                             driver.find_element_by_css_selector('...').click()
                             i += 1
               with open('C:/....txt','w',encoding='UTF-8') as file:
                      for j in text1:
                             file.write(j+'\n')
```

4月 3日~4月 5日

I - 4. 데이터



nam_yeeeeh갤럭시s9+죠흐당♡♥

(5) (4)

#컵을 #커플폰 #기념일 #3주년 #갤럭시s9플러스 #아이폰언제써보냐 #듀얼카메라 #화질고와 #기념셀카 #귀요미 #스노우카메라 #봄 #라이더자켓 #영화 #곤지암 #제발좀가지말라고 #말좀들어라 #존나무서워 #가장맛있는족발 #육회공작소 #마이쪙 #먹스타그램 #데일리 #일상 #데이트 #좋아요 #이디야 #당보충 #바닐라라떼 umcertain장조림과 수육사이 차슈. 문동이 일본가서 사온 라멘. 반숙을 원했지만 완숙이 된 계란. 드디어 사용한 스티 치 접시. 이디야 치즈스틱케익이 없어서 초코 구입. 고양이 젓가락받침대와 많이 세탁해서 바랜 밤색 린넨. 집에는 버스커버스커 재생중.주말오전 the_court_clerk8#20180323 아침부터 혼자 바쁨돋았음 ㅎㅎ 동전지갑 득템하고요(이어폰 케이스로 찰떡)



#컵흘 #커플폰 #기념일 #3주년 #갤럭시s9플러스 #아이폰언제써보냐 #듀얼카메라 #화질고와 #기념셀카 #귀요미 #스노우카메라 #봄 #라이더자켓 #영화 #곤지암 #제발좀가지말라고 #말좀들어라 #존나무서워 #가장맛있는족발 #육회공작소 #마이쪙 #먹스타그램 #데일리 #일상 #데이트 #좋아요 #이디야 #당보충 #바닐라라떼

Ⅱ. 데이터 정제



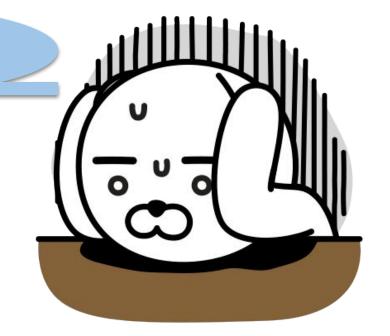
#인스타그램 #데일리 #일상 #daily #dailylook #ootd #instadaily #selfie #selca #like4like #셀 스타그램 #얼스타그램 #셀카 #셀피 #카페 #cafe #카페스타그램 #이디야

#이디야카페 #이디야커피 #카페스타그램 #이 디야 #손그림 #그림 #art #doodle #penart #drawing #ediyacoffee #ediya #followme #like4like #좋아요반사 #followforfollow #선팔 하면맞팔 #언팔※#낙서

#99 #00 #20 #19 #고딩 #여고생 # 수원역 #인친 #데일리 #일상 #소통 #셀스타그 램 #selfie #like4like #likeforlike #followme #follow #고3 #셀기꾼 #팔로우 #팔로미 #모모 스테이크 #이디야 #친구 #일상 #거울샷 #폭풍업뎃 #셀스 타그램 #이디야 #카페 #시험기간 #고딩 #좋아 요반사 #좋아요 #맞팔 #소통 #selfie #cafe #mirrorselfie #|4| #f4f #like4like

"제발 통일 좀 해줘..."

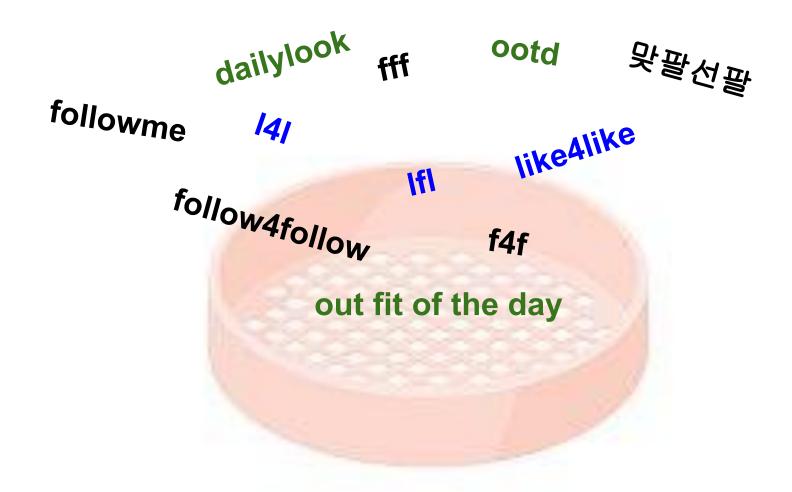
#일상 #일상그램 #데일리 #데일 리그램 #데일리룩 #데일리라이프 #라이프 #라 이프그램 #daily #dailygram #dailypic #dailylook #dailylife #dailylifegram #starbucks #스타벅스 #스벅md #스벅텀블러 #스벅 #스 타벅스투톤핑크사이렌보온병 #그러데이션사 이렌핸들글라스 #그라데이션사이렌핸들글라

























STARBUCKS



EDIYA



III. WORD CLOUD 2



STARBUCKS



EDIYA





공통적으로 사진, 셀피多

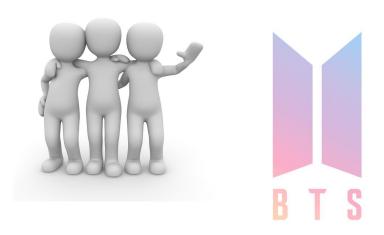
스타벅스

친구

방탄소년단

이디야

육아 연애중





C

STARBUCKS & EDIYA







STARBUCKS



BTS Score Top 10 Debut on Billboard 200 With 'Love Yourself: Her' Album

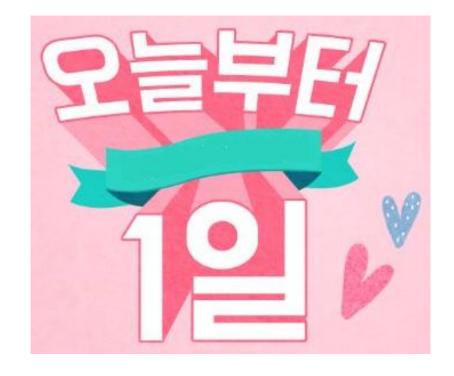




EDIYA

이구동성 리즈 이벤트





& THANK YOU

