



DIGITAL  
VINYL

DESIGNED  
BY  
L@RGO  
ADSTORE

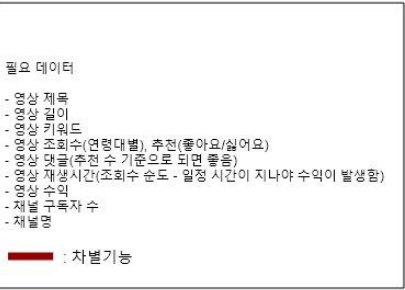
An abstract graphic on the left side of the slide. It features a large circle divided horizontally into a blue top half and an orange bottom half. The word "INDEX" is written in white capital letters across the center of this circle. Surrounding the circle are several concentric white circles. In the background, there are white triangles of various sizes and orientations, some pointing up and some pointing down, scattered across a gradient of blue and orange. The overall design is modern and geometric.

INDEX

-  
목적

-  
웹 크롤링

-  
분석 결과



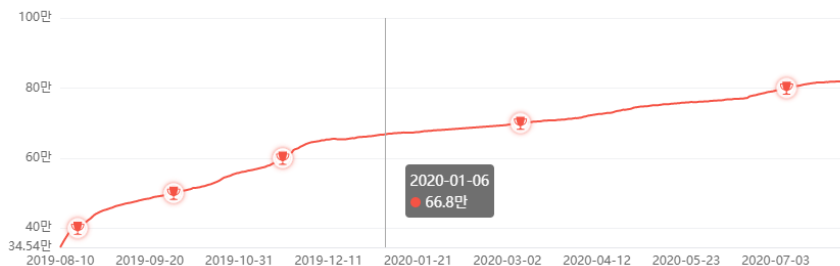


# 채널 정보 웹 크롤링

구독자 히스토리 데이터 (최근 1년)

일별 데이터    누적 데이터

안내: 최신 YouTube 공개 구독자 수 표시 방식에 따라 앞 3자리만 표시되기 때문에 구독자 수 그래프 변화가 있습니다



```
# 구독자 누적 차트에 마우스 호버 실행
element_to_hover_over = driver.find_element_by_xpath('//*[@id="channel-history-sub-chart"]/div[1]/canvas')
hover = ActionChains(driver).move_to_element(element_to_hover_over) # 차트의 가운데로 마우스 호버
hover.perform()

# 차트 사이즈 가져오기
element_size = element_to_hover_over.size
element_height = element_size['height']
element_width = element_size['width']

# 차트의 좌측으로 이동
x_offset = (element_width / 2)
move = ActionChains(driver).move_by_offset(-x_offset, 0)
move.perform()

# 차트의 우측으로 이동하면서 날짜/누적 구독자수 저장
'''
soup = BeautifulSoup(html, 'lxml') -> 사용해서 크롤링 시간 단축 해보기!
'''

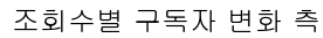
list_sub_date = []
list_subscribe_cnt = []
for idx in range(int(element_width / 2)):
    # 마우스 이동
    move = ActionChains(driver).move_by_offset(2, 0)
    move.perform()

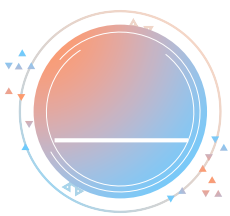
    chart_info = driver.find_element_by_xpath('//*[@id="channel-history-sub-chart"]/div[2]').text
    if chart_info.find('\n') != -1:
        list_sub_date.append(chart_info.split('\n')[0])
```



# 웹 크롤링 결과

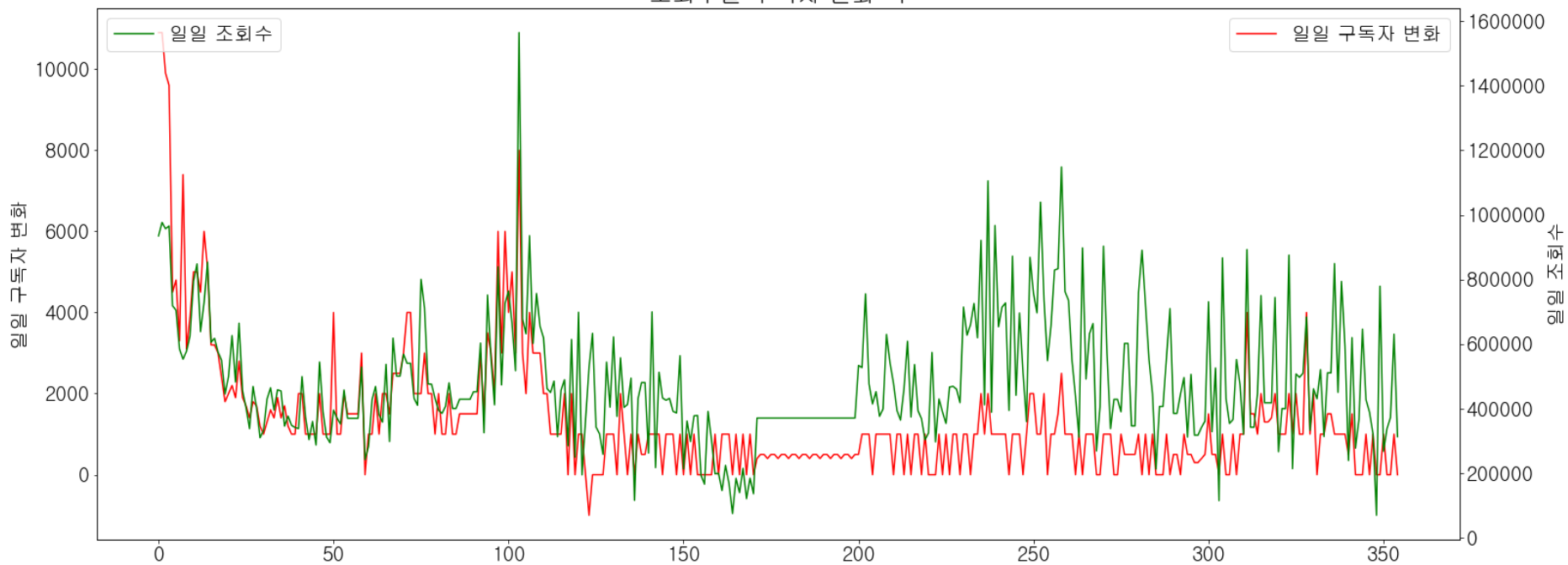
A	B	C	D	E
	date	subscribe count	daily subscribe count	daily view count
0	2019-08-06	168700	1000	353200
1	2019-08-07	169700	1300	423900
2	2019-08-08	171000	1200	399300
3	2019-08-09	172200	1400	442300
4	2019-08-10	173600	700	279500
5	2019-08-11	174300	800	286200
6	2019-08-12	175100	1000	527200
7	2019-08-13	176100	800	404500
8	2019-08-14	176900	400	243200
9	2019-08-15	177300	500	296800
10	2019-08-16	177800	400	285700
11	2019-08-17	178200	400	326400
12	2019-08-18	178600	400	366700
13	2019-08-19	179000	500	347100





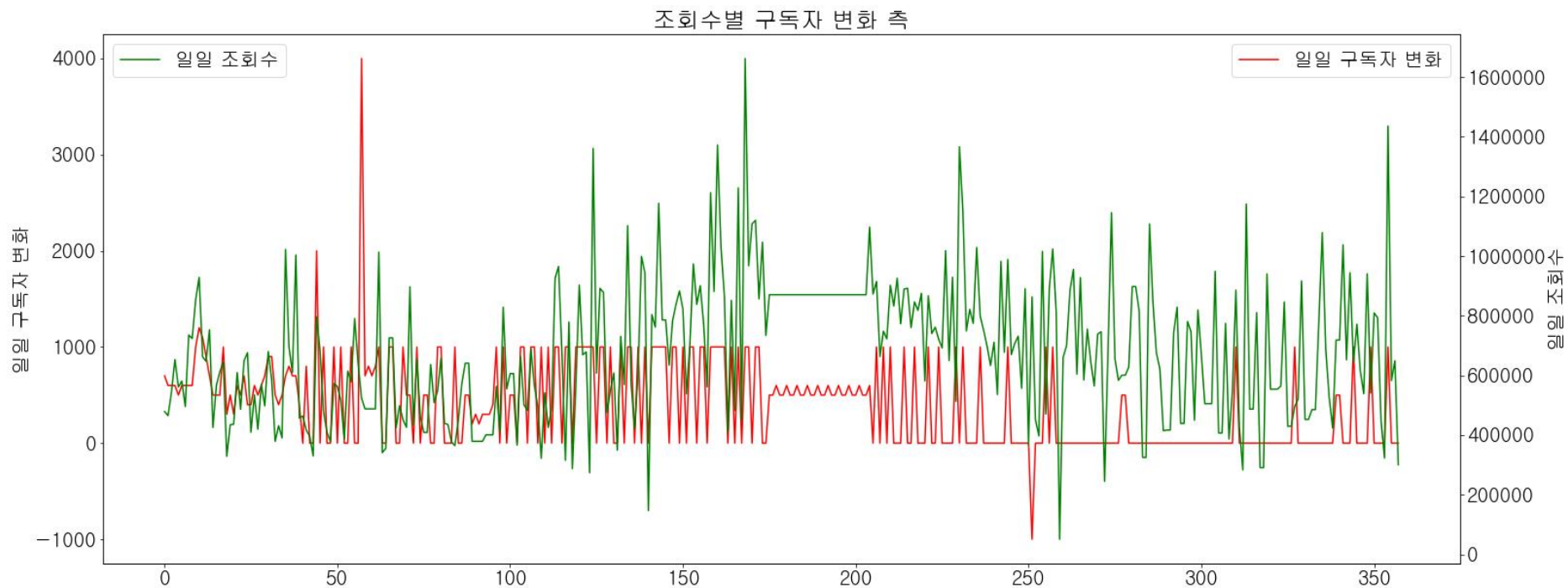
# 오 킹 TV 채널

조회수별 구독자 변화 측

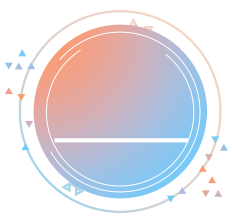




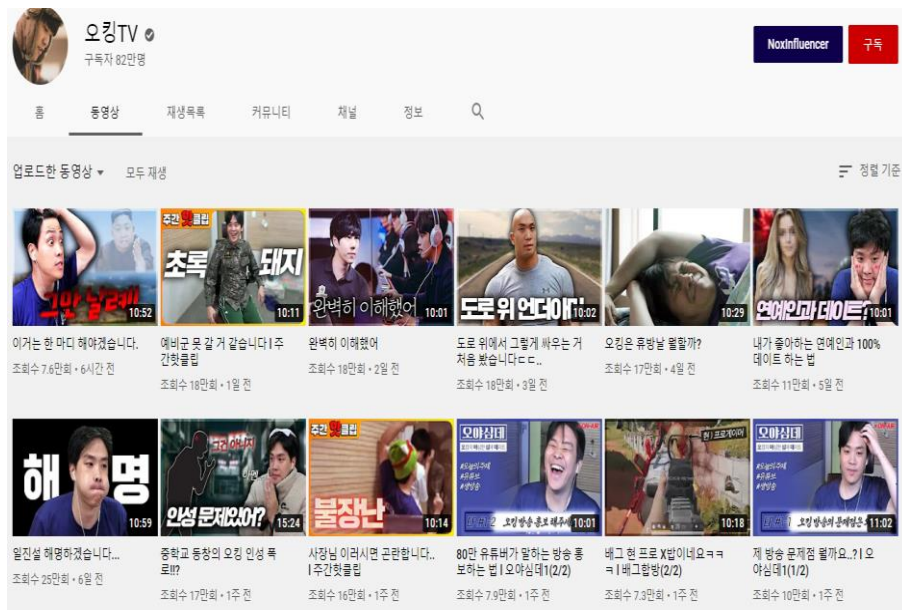
# 두찌와뿌꾸 채널







# 채널 업로드 정보 웹크롤링



```
# 동영상 페이지 제일 밑으로 스크롤
while True:
    last_height = driver.execute_script('return document.documentElement.scrollHeight')
    # 현재 화면의 길이를 리턴 받아 last_height에 넣음
    for i in range(10):
        body.send_keys(Keys.END)
        # body 본문에 END키를 입력(스크롤내림)
        time.sleep(SCROLL_PAUSE_TIME)
    new_height = driver.execute_script('return document.documentElement.scrollHeight')
    if new_height == last_height:
        break;

page = driver.page_source
soup = BeautifulSoup(page, 'lxml')

# 채널명 크롤링
username = soup.find('div', {'id': 'text-container'}).text
username = username.strip()

# 제목, url 크롤링
all_videos = soup.find_all(id='dismissable')

list_title = []
list_url = []
list_video_length = []
for video in all_videos:
    title = video.find(id='video-title')
    if len(title.text.strip()) > 0: # 공백을 제거하고 글자수가 0보다 크면 append
        list_title.append(title.text)

    #find('a', {'id': 'thumbnail'})['href']
    url = video.find(id='video-title')['href'] # url append
    list_url.append(url)

    video_lenth = video.find('span', {'class': 'style-scope ytd-thumbnail-overlay-time-status-renderer'})
    list_video_length.append(video_lenth.text.strip())
```

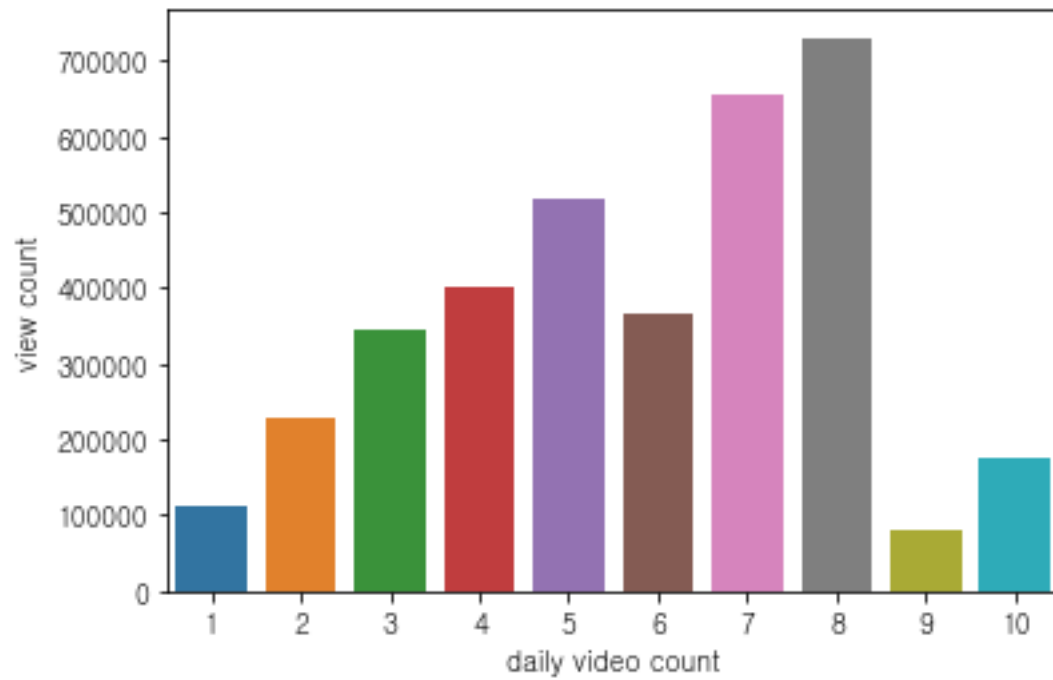


# 웹 크롤링 결과

	A	B
1	upload_date	daily video count
2	2017-12-09	1
3	2017-12-19	1
4	2017-12-30	1
5	2018-01-01	1
6	2018-01-04	1
7	2018-01-07	1
8	2018-01-11	1
9	2018-01-14	1
10	2018-02-13	1
11	2018-02-18	1
12	2018-02-25	1
13	2018-03-02	1
14	2018-03-15	1
15	2018-03-17	1
16	2018-03-22	1
17	2018-03-30	1
18	2018-04-03	1
19	2018-04-07	1

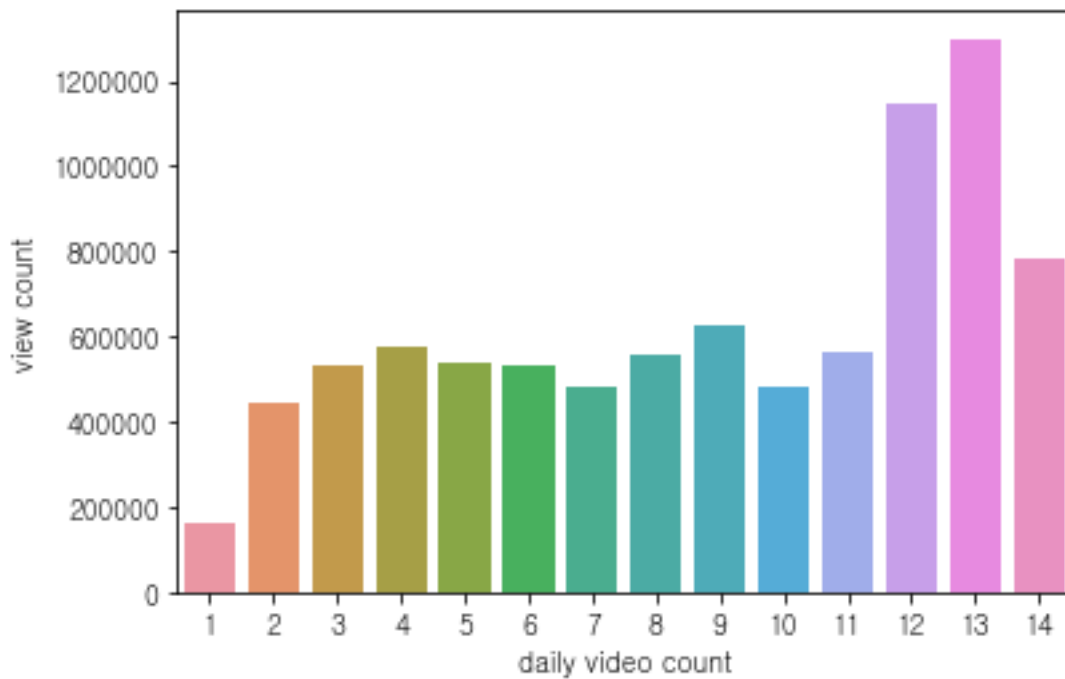


# 이 스타 TV 채널 영상 개 수 별 조 회 수 시 각 화





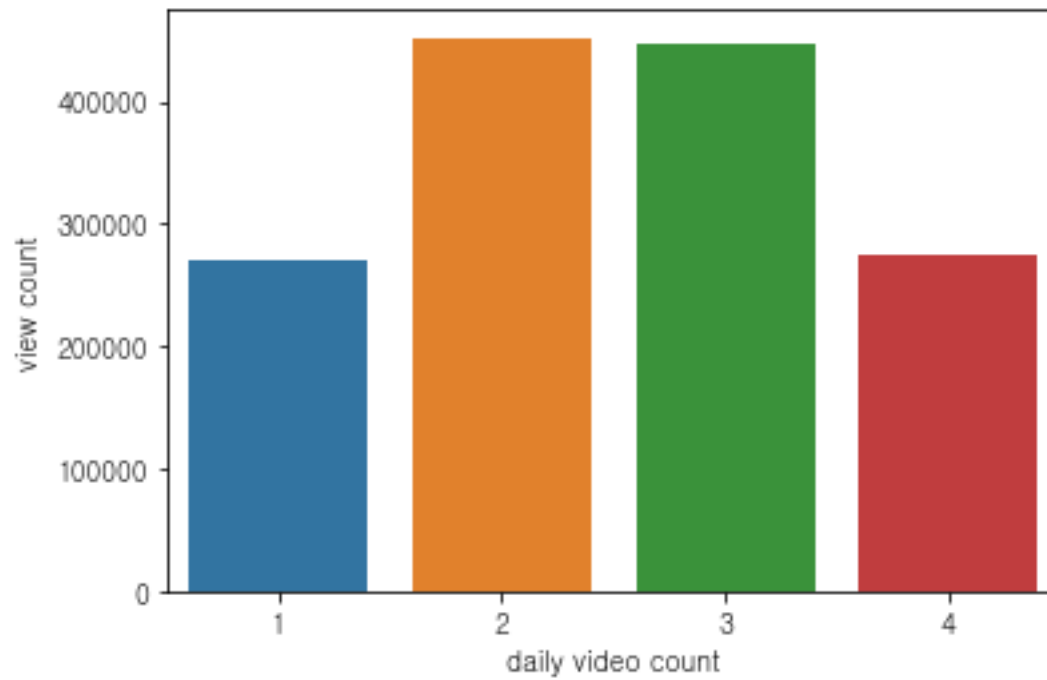
# 두 씨 와 뿌 구 채널 수 시각화





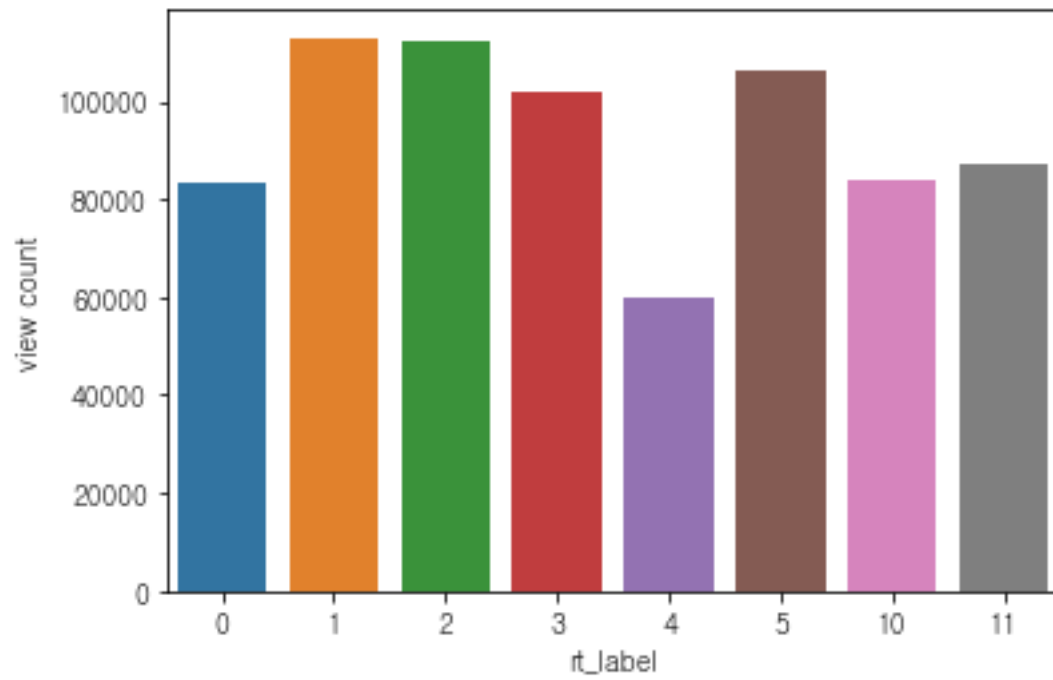
오 킹 TV 채널  
영상 개 수 별

조 회 수 시 각 화



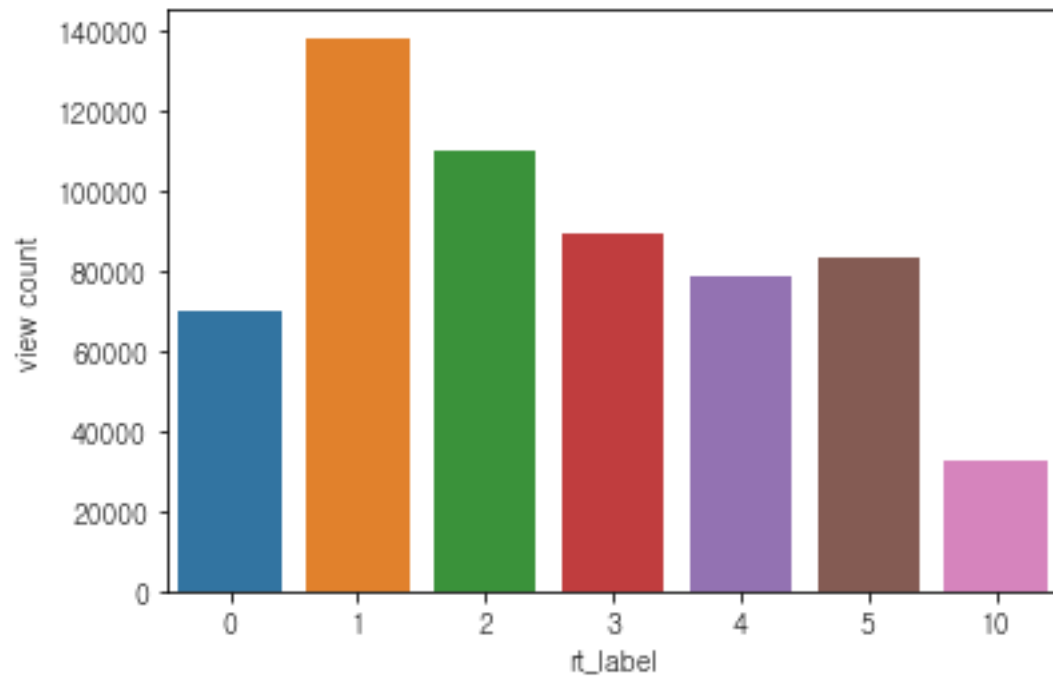


# 이 스타TV 채널 영상 길이 별 조회 수 시각화





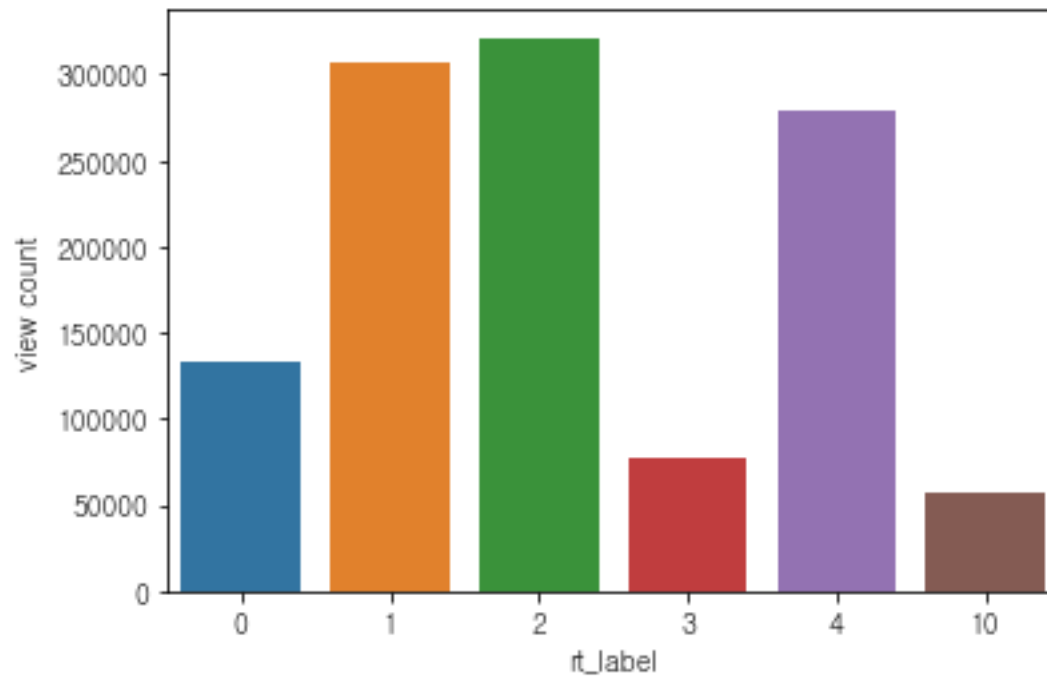
# 두 씨와 뿌꾸 채널 영상 길이 별 조회 수 시각화





오 킹 TV 채널  
영상 길이 별

조회 수 시각화







# 누적 구독자 변화 예측 시각화

```
subscribe = list(df_read_channel_info['subscribe count'].values) # 누적 구독자 데이터 추출
seq_len = 50 # 기준 날짜 설정(50 -> 최근 50일)
# 데이터 분리 및 정규화
X, Y = normalized_df(subscribe, seq_len)

# 데이터 분할
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X,
                                                    Y,
                                                    test_size = 0.1,
                                                    random_state = 0)
X_train = np.reshape(X_train, (X_train.shape[0], X_train.shape[1], 1))
X_test = np.reshape(X_test, (X_test.shape[0], X_test.shape[1], 1))

# 모델 구성
model = Sequential()
model.add(LSTM(100, return_sequences=True, input_shape=(seq_len, 1))) # 왜 input_shape를 2
model.add(LSTM(50, return_sequences=False))
model.add(Dense(1, activation='linear'))
model.compile(loss='mse', optimizer='adam')

model.fit(X_train, Y_train, validation_data=(X_test, Y_test), batch_size=1, epochs=10)

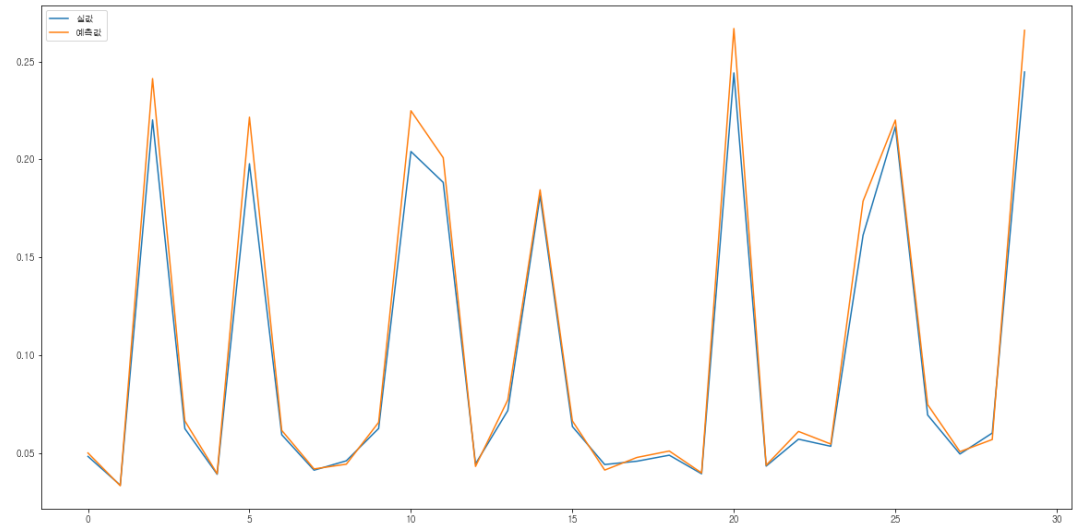
# 예측 값과 실제 값의 비교
Y_prediction = model.predict(X_test).flatten()
for i in range(10):
    label = Y_test[i]
    prediction = Y_prediction[i]
    print("실제가격: {:.3f}, 예상가격: {:.3f}".format(label, prediction))

from sklearn.metrics import r2_score
print('결정계수', r2_score(Y_test, Y_prediction))
```



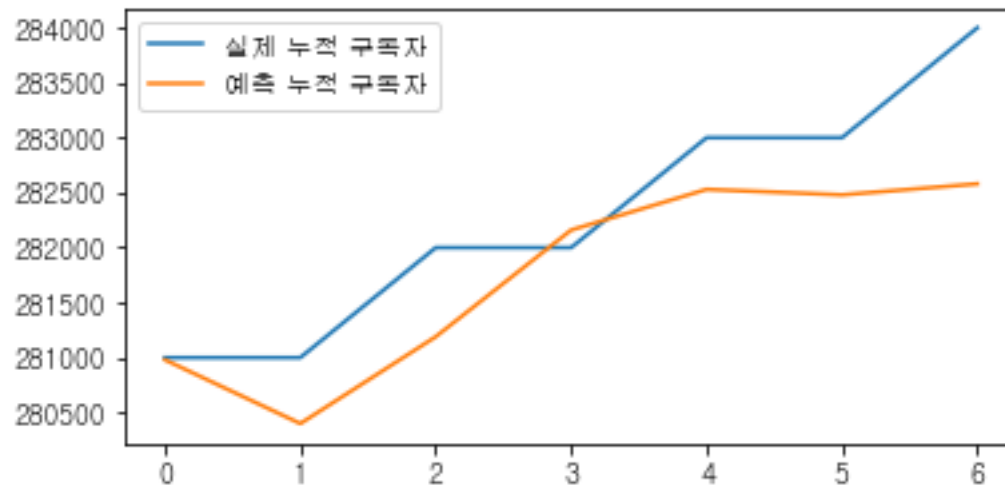
# 누적 구독자 변화 예측 시각화

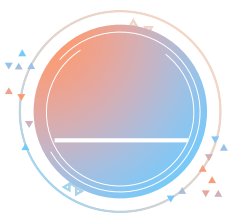
```
val loss: 1.0165e-04  
실제가격: 0.048, 예상가격: 0.050  
실제가격: 0.034, 예상가격: 0.033  
실제가격: 0.220, 예상가격: 0.241  
실제가격: 0.062, 예상가격: 0.066  
실제가격: 0.039, 예상가격: 0.039  
실제가격: 0.198, 예상가격: 0.222  
실제가격: 0.059, 예상가격: 0.062  
실제가격: 0.041, 예상가격: 0.042  
실제가격: 0.046, 예상가격: 0.044  
실제가격: 0.062, 예상가격: 0.066  
결정계수 0.9809205139992156
```



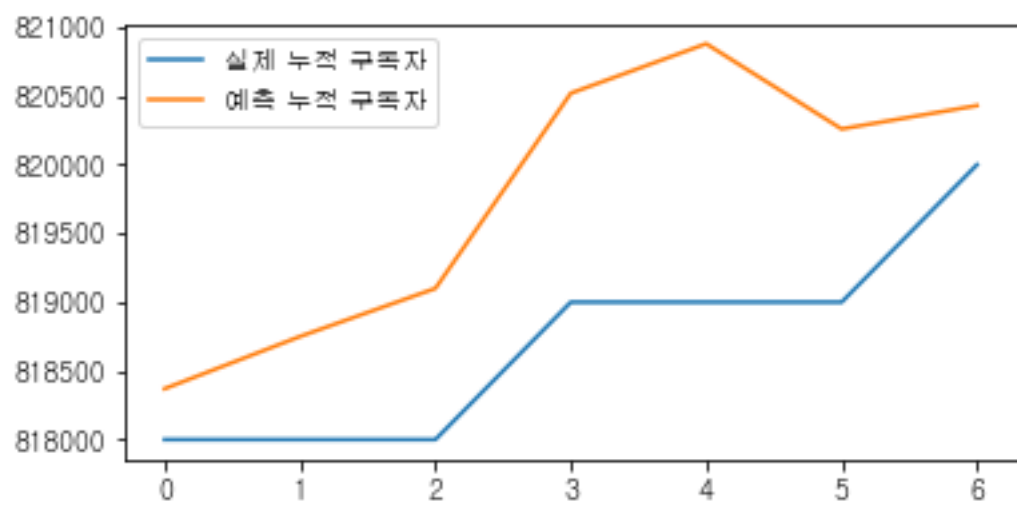


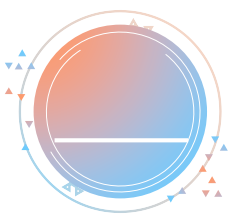
# 이 스타 TV 채널 영상 개 수 별 조 회 수 시 각 화



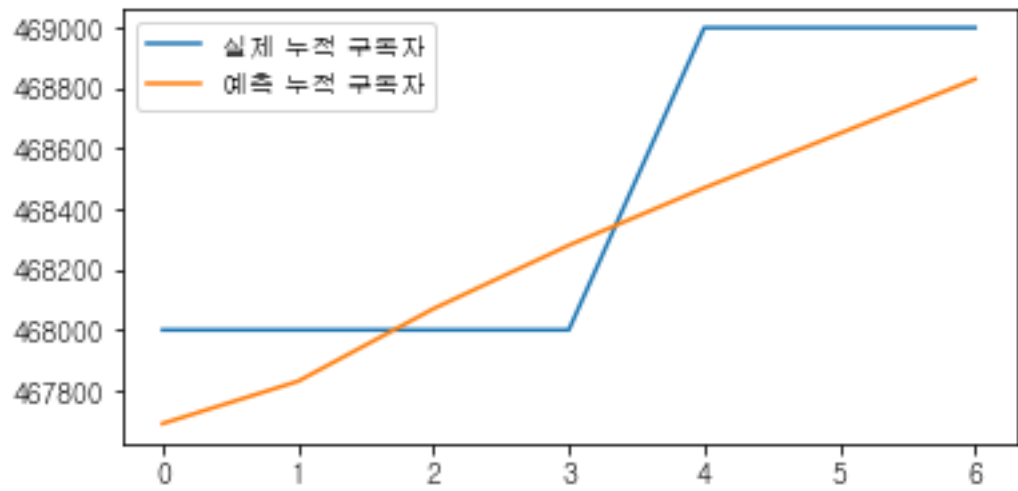


# 오 킹 TV 채널 누적 구독자 변화 예측



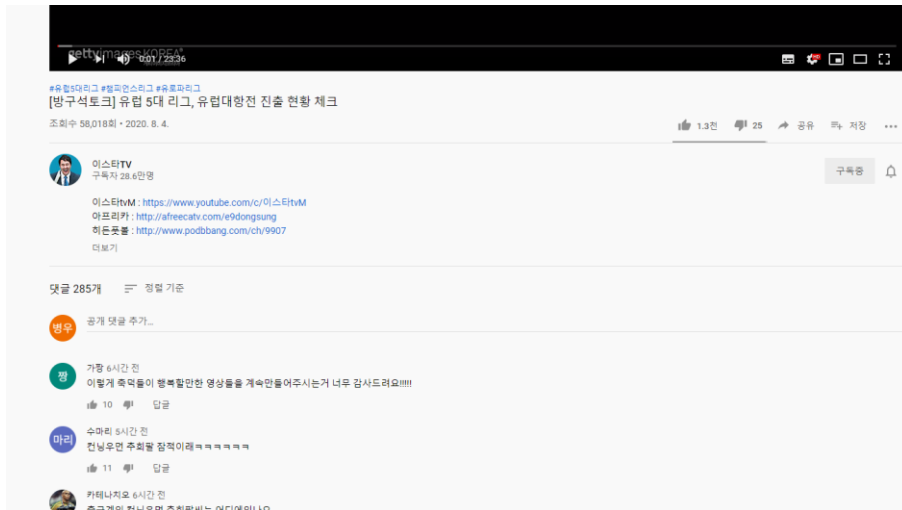


# 두 씨 와 뿌 꾸 채 널 누 적 구 독 자 변 화 예 측





# 댓글 웹크롤링



```
SCROLL_PAUSE_TIME = 0.5 # 한번 스크롤 하고 멈출 시간 설정
# 더보기 버튼 클릭을 위해 최초 1회 스크롤 다운
body.send_keys(Keys.PAGE_DOWN)
time.sleep(SCROLL_PAUSE_TIME)

# 더보기 클릭
more_button = driver.find_element_by_xpath('//*[@id="more"]/yt-formatted-string')
more_button.click()

# 크롤링을 위해 화면 맨 아래까지 스크롤 내리기
while True:
    last_height = driver.execute_script('return document.documentElement.scrollHeight')
    # 현재 화면의 길이를 리턴 받아 last_height에 넣음
    for i in range(10):
        body.send_keys(Keys.PAGE_DOWN)
        # body 본문에 END키를 입력(스크롤내림)
        time.sleep(SCROLL_PAUSE_TIME)
        new_height = driver.execute_script('return document.documentElement.scrollHeight')
        if new_height == last_height:
            break;

html = driver.page_source
soup = BeautifulSoup(html, 'lxml')

# 댓글 id 추출
all_comment_id = soup.find_all('a', {'id' : 'author-text'})
comment_id = [soup.find_all('a', {'id' : 'author-text'})[n].text for n in range(0, len(all_comment_id))]
comment_id = [i.strip() for i in comment_id]

# 댓글 추출 (답글 제외)
all_comment_contents = soup.find_all('yt-formatted-string', {'id' : 'content-text'})
comment_contents = [soup.find_all('yt-formatted-string', {'id' : 'content-text'})[n].text for n in range(len(all_comment_contents))]
comment_contents = [i.strip() for i in comment_contents]
```



# 웹 크롤링 결과

A	B	
	comment id	
0	milk 0423	섭종각이라던데..형 이제라도 돈아끼고 키티누님이랑 데이트하는영상좀올려줘
1	...	5:55 농 호나우두 금카
2	최준우	1:25 홀리~~ 췌!!4:29 홀리~~ 췌!!
3	조정현	5:55 채팅 예언보소
4	이름	5:00실수하네
5	오.	6:18 ㅅㅂㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋ
6	으으음	뿌꾸형 거만흑우일때 넘나 기여워욥ㅋㅋ 개육겨 증말 홀리췌~! 홀리 췌~! 개재밋서
7	차재근	5:48 믹슈 손가락 까닥까닥하는거 호나우두 셀레브레이션 복선 지렸다,,,ㄷㄷ
8	환	피자먹방 보고 온사람
9	김덕현	2:28 날강두 그 발언..
10	노해석	와 호나우두 8카가 나오긴 하는구나...난 EP 500만원 받고 끝인데...
11	유영빈	2:29 실수하는 소리하네
12	김민수	8:08 TB음바페 8:25 라이브음바페진짜 어제 생방으로 봤는데 쌍바페 레전드다 ㅋㅋ
13	추배	믹슈형 아이콘 포함 88팩에서도 아이콘 띄우고호나우두 금카도 띄우고 금손인가?



# 이모티콘, 불용어 처리 전

df_read_comment - DataFrame			
Index	ame	comment id	comment
0	0	이스타TV	제모크림 구매 링크 : <a href="http://s.godo.kr/lzgk">http://s.godo.kr/lzgk</a>
1	1	새참	첼시는 뭐 갈락티코도 아니고 존나 공격쪽만 보강하네 ㅋㅋㅋㅋ
2	2	khm	케파를 이스타가 감 > 케파 안정적임 > 케파 잘하는거보고 이...
3	3	이쁜우리 누님	이야 첼시팬들은 우동사리 나가고 나서, 완전 기쁨 그자체네요. 축하드립니다
4	4	차감홍	
5	5	호호호	서로 미끄러지... 감미네 ㅋㅋ 고려의 호호하 스페는 스토하
6	6	anonymous	이 영상을 장지현해설이 좋아합니다
7	7	김동형	첼시팬은 행복합니다😊
8	8	이현수	ㅋㅋㅋ 공격 보다가 수비보니까 무게감 많이 떨어지네
9	9	탁원준	이게 될거라 생각 못했는데 진짜 되면 대박이다
10	10	[Hari st...	하베르츠까지 지르면 첼시 진짜 p1 너무 재밌어지겠네요 ㄷㄷ...
11	11	봉 가	뤼디거 장사 잘한다.국대 친구들 모으네
12	12	파드렘	첼시 챔스 진출성공시 정확한 모르겠지만 800억~900억의 이적 자 금을 확보 할 수 있습니다
13	13	Supreme love	첼시 미쳤네.. 0입 한시즌 했다고 이렇게 폭풍영입하나..리버풀도 미나미노만... 0.02 첼시팬이데 고려지 리스F보고 11도 그렇게 인기라 서로서르





# 이모티콘, 불용어 처리 후

comment\_result - List (352 elements)

Indx	Type	Size	Value
0	str	1	제모크림 구매 링크 : <a href="http://sgodokr/lzgk">http://sgodokr/lzgk</a>
1	str	1	헬시는 뭐 갈락티코도 아니고 존나 공격쪽만 보강하네
2	str	1	케파를 미스타가 꺾 > 케파 안정적임 > 케파 잘하는거보고 미스타가 사과할 > 케파가 그뒤로 못함 결론: 이 ...
3	str	1	이야 헬시팬들은 우동사리 나가고 나서 완전 기쁨 그자체네요 축하드립니다
4	str	1	마 아스날 한번 와라 어 어디서 좋은 일만 할려고 젊을 때 고생해야지
5	str	1	옛날 마드리드 느낌나네공격은 초호화 수비는 초토화
6	str	1	이 영상을 장지현해설이 좋아합니다
7	str	1	헬시팬은 행복합니다
8	str	1	공격 보다가 수비보니까 무게감 많이 떨어지네
9	str	1	이게 될거라 생각 못했는데 진짜 되면 대박이다
10	str	1	하베르츠까지 지르면 헬시 진짜 p1 너무 재밌어지겠네요

Save and Close Close



# 명사 추출

```
def get_noun(comment_txt):  
    twitter = Twitter()  
    noun = []  
  
    if len(comment_txt)>0:  
        tw = twitter.pos(comment_txt)  
        for i,j in tw:  
            if j == 'Noun':  
                noun.append(i)  
    return noun
```

df_comment_result - DataFrame				
Index	mm	token		
0	제...	['제모', '크림', '구매', '링크']		
1	첼...	['첼시', '뮌', '갈락티코', '존나', '공격', '족', '보강']		
2	케...	['케파', '이스', '타가', '감', '케파', '안정', '적임', '...		
3	미...	['첼시', '팬', '우동', '사리', '완전', '기쁨', '자체', '축하']		
4	마 ...	['마', '아스날', '한번', '일만', '때', '고생']		
5	옛...	['옛날', '마드리드', '느낌', '공격', '초', '호화', '수비', '초토화']		
6	미 ...	['미', '영상', '장지현', '해설']		
7	첼...	['첼시', '팬']		
8	공...	['공격', '비보', '무게', '감']		
9	미...	['미', '생각', '진짜', '대박']		
10	하...	['하베르츠', '첼시', '진짜']		
11	뤼...	['뤼', '디거', '장사', '국', '친구']		
12	첼...	['첼시', '첼스', '진출', '공시', '정확', '힌', '미적', '자금', '확보', '수']		
13	첼...	['첼시', '미쳤', '입', '시즌', '했다', '폭풍', '입하', '리버풀', '미나', '미노', '입했', '디거', '입', '왜']		

Format    Resize    Background col    Column min/max    Save and Close    Close



# 명사 빈도수 확인

```
noun_list = []
for i in range(len(df_comment_result)):
    for j in range(len(df_comment_result['token'].iloc[i])):
        noun_list.append(df_comment_result['token'].iloc[i][j])

counts = Counter(noun_list) # 추출된 명사 빈도수 확인
tags = counts.most_common(30) # 빈도수 상위 30개 추출
###
```

tags - List (30 elements)

Indx	Type	Size	Value
0	tuple	2	('헬시', 122)
1	tuple	2	('하베르츠', 68)
2	tuple	2	('수비', 58)
3	tuple	2	('입', 35)
4	tuple	2	('진짜', 30)
5	tuple	2	('케파', 27)
6	tuple	2	('팬', 23)
7	tuple	2	('선수', 23)
8	tuple	2	('골', 21)
9	tuple	2	('시즌', 20)
10	tuple	2	('베르너', 19)

Save and Close Close



워드 클라우드

수비수 마운트 하  
사 「체 시」 격진  
스타 골 좀더 세아 배  
입 진짜 케파 추 이적  
팬  
하 베 르 후  
제발 맨유 센터백 다음 팀

# Question

