## 회사 밖의 데이터 - 데이터로 사이드 프로젝트 만들기 나의 브런치 구독자 데이터 분석기

파이썬을 활용한 데이터 분석기

강원양



# 회사 안의 데이터

01 하는 일



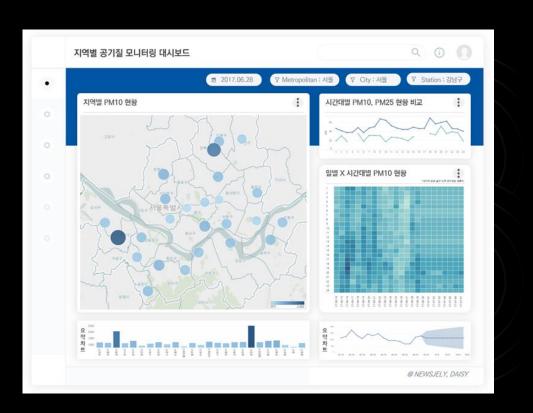
데이터 시각화를 기획하고 만들고 이야기 합니다.



데이터 **시각화를 기획**하고 만들고 이야기 합니다.



동일한 데이터도 활용 목적에 따라 더 효과적인 방법으로 표현

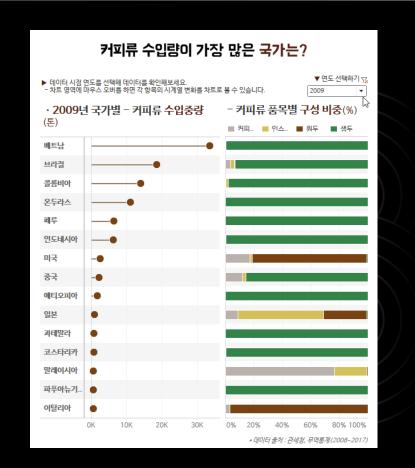


데이터 시각화를 기획하고 만들고 이야기 합니다.



하는 일

시각화를 통한 데이터 분석
시각적 분석 기반 인사이트 도출
데이터 스토리텔링 콘텐츠



## 회사 밖의 데이터

02 프로젝트 배경



**SOSCON 2019** 

SAMSUNG OPEN SOURCE CONFERENCE 2019

#### 왜 프로젝트를 하게 되었나?

- 데이터를 활용해 궁금증을 해결하기
- 현실적인 한계를 극복하기 위한 시도하기
- 데이터 분석 과정에서 시각화의 역할 이야기하기



#### 어떻게 프로젝트를 하게 되었나?

- 파이썬 6주 수업
- 선생님과 함께하는 미니 프로젝트

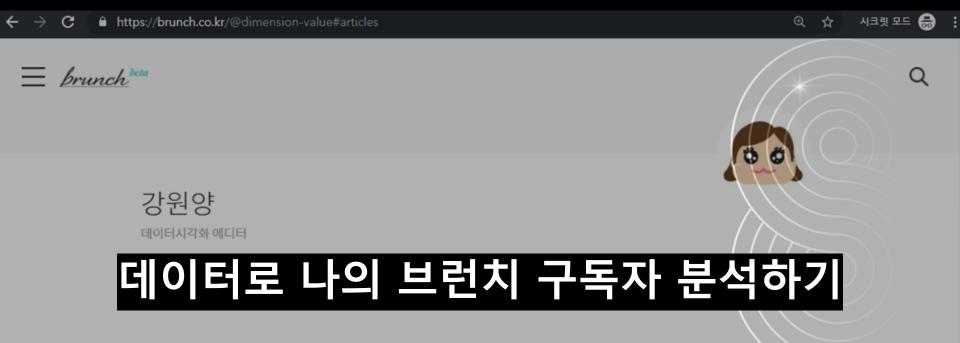


## 회사 밖의 데이터

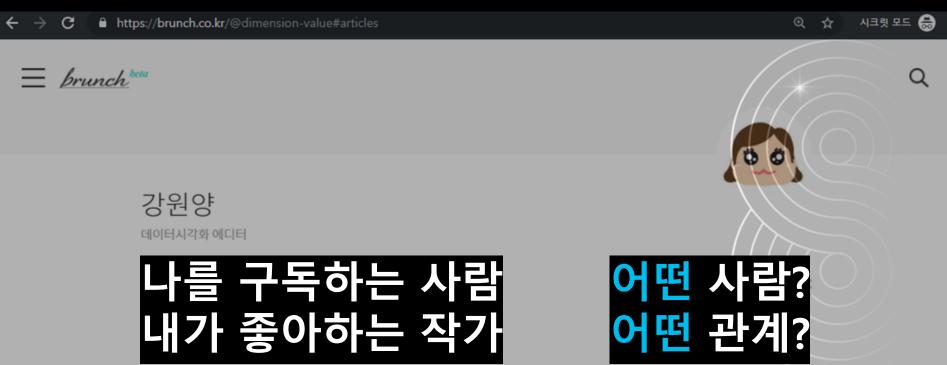
03 프로젝트 개요



#### 프로젝트 개요



프로젝트 개요



#### 프로젝트 개요



+ableau<sup>†</sup>†public



## 회사 밖의 데이터

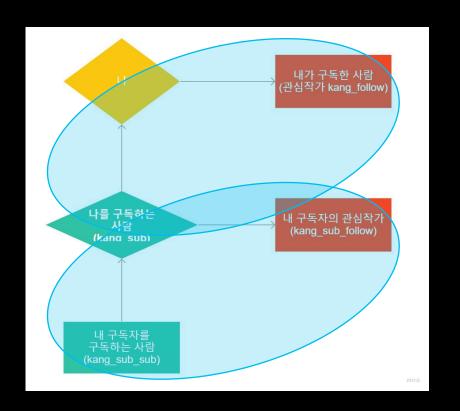
03 프로젝트 내용



**SOSCON 2019** 

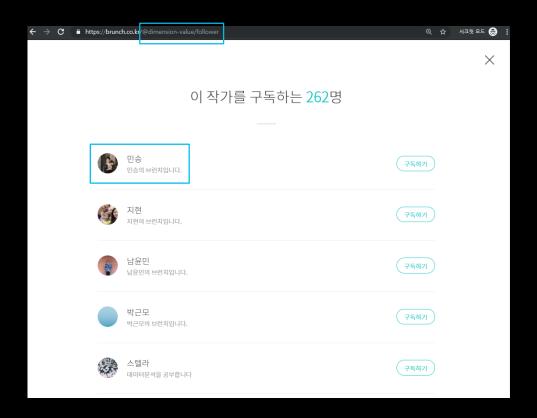
SAMSUNG OPEN SOURCE CONFERENCE 2019

#### 파이썬을 활용한 데이터 수집 및 정제





#### 웹 크롤링 : 구독자/ 관심작가 리스트





#### 웹 크롤링 : 구독자/ 관심작가 리스트

```
1 from selenium import webdriver
            from selenium.webdriver.common.kevs import Keys
            from bs4 import BeautifulSoup
          4 import time
         1 #구독자는 follower, 관심작가는 following
In [2]:
            sub_url = "https://brunch.co.kr/@dimension-value/follower
        | driver = webdriver.Chrome('./drivers/chromedriver.exe')
         2 driver.get(sub_url)
In [4]: 1 main = driver.find_element_by_tag_name("body")
In [5]:
          1 #무한 스크롤링
            num_of_pagedowns =50
            while num_of_pagedowns:
                main.send_keys(Keys.PAGE_DOWN)
                time.sleep(0.3)
                num_of_pagedowns -=1
        1 html = driver.page_source
         1 | soup = BeautifulSoup(html, 'html.parser')
        | | profiles = soup.select('ul.list_follow > li')
         2 print(type(profiles), len(profiles))
        <class 'list'> 259
        1 index = 0
            sub_data=[]
           for sub_profile in profiles:
               sub_name = soup.select('.tit_subject')[index].text
sub_desc = soup.select('.desc_follow')[index].text
                sub_link=soup.select('.link_follow')[index]['href']
               print('/@dimension-value', index, sub_name, sub_desc.strip(), sub_link, sep = ' | ')
                sub_data.append(['/@dimension-value', index, sub_name, sub_desc.strip(), sub_link,])
        /@dimension-value | 0 | 스텔라 | 데이터분석을 공부합니다 | /@ivoryohO609
        /@dimension-value | 1 | castellan | castellan의 브런치입니다. | /@hancastellan
        /@dimension-value | 2 | KyeongMi | KyeongMi의 브런치입니다. | /@kyeongmimfu8
        /@dimension-value | 3 | 양인환 | 양인환의 브런치입니다. | /@reay1212
        /@dimension-value | 4 | 박구라 | 배우는 법을 배워가는 중입니다. | /@gurapark
```



#### 웹 크롤링 : 구독자/ 관심작가 리스트

```
from selenium import webdriver
            from selenium.webdriver.common.kevs import Keys
            from bs4 import BeautifulSoup
          4 import time
In [2]:
         1 #구독자는 tollower, 관심작가는 tollowing
            sub_url = "https://brunch.co.kr/@dimension-value/follower"
        1 driver = webdriver.Chrome('./drivers/chromedriver.exe')
         2 driver.get(sub_url)
In [4]: 1 main = driver.find element by tag name("hody")
In [5]:
           #무한 스크롤링
            num_of_pagedowns =50
            while num_of_pagedowns:
               main.send_kevs(Kevs.PAGE_DOWN)
               time.sleep(0.3)
               num_of_pagedowns -=1
         1 html = driver.page_source
         1 soup = BeautifulSoup(html, 'html.parser')
        | | profiles = soup.select('ul.list_follow > li')
         2 print(type(profiles), len(profiles))
        <class 'list' > 259
        1 index = 0
            sub_data=[]
            for sub_profile in profiles:
               sub_name = soup.select('.tit_subject')[index].text
               sub_desc = soup.select('.desc_follow')[index].text
               sub_link=soup.select('.link_follow')[index]['href']
               print('/@dimension-value', index, sub_name, sub_desc.strip(), sub_
               sub_data.append(['/@dimension-value', index, sub_name, sub_desc.st
               index = index + 1
       /@dimension-value | 0 | 스텔라 | 데이터분석을 공부합니다 | /@ivoryoh0609
       /@dimension-value | 1 | castellan | castellan의 브런치입니다. | /@hancastel
       /@dimension-value | 2 | KyeongMi | KyeongMi의 브런치입니다. | /@kyeongmimfu
       /@dimension-value | 3 | 양인환 | 양인환의 브런치입니다. | /@reay1212
        /@dimension-value | 4 | 박구라 | 배우는 법을 배워가는 중입니다. | /@gurapar
```

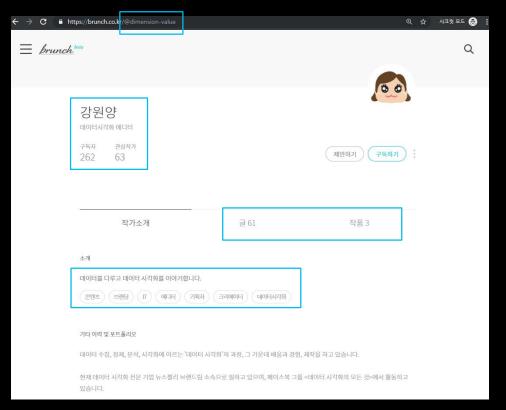
```
In [36]
             driver = webdriver.Chrome('./drivers/chromedriver.exe')
             total_sub_data=[]
             for i in sub_url
                 driver, get (i)
                 time.sleep(3)
                 SCROLL_PAUSE_TIME = 1.5
                 # Get scroll height
                 last_height = driver,execute_script("return document,body,scrollHeight")
         14
                 while True:
                     # Scroll down to bottom
         16
                     driver, execute_script("window, scrollTo(0, document, body, scrollHeight);")
         18
                     # Wait to load page
                     time,sleep(SCROLL_PAUSE_TIME)
                     # Calculate new scroll height and compare with last scroll height
                     new_height = driver,execute_script("return document,body,scrollHeight")
                     if new_height == last_height:
         24
                         break
         25
                     last_height = new_height
         26
                 html = driver,page_source
         28
                 soup = BeautifulSoup(html, 'html,parser')
         29
         30
                 profiles = soup.select('ul.list_follow > li')
         31
         32
                 index = 0
                 for sub_profile in profiles:
         34
                         sub_name = soup,select(',tit_subject')[index],text
         35
                         sub_desc = soup,select(',desc_follow')[index],text
         36
                         sub_link=soup.select(',link_follow')[index]['href']
         37
                         print('/'+i,split('/')[3], index, sub_name, sub_desc,strip(), sub_link, sep = ' | ')
         38
                         total_sub_data,append(['/'+i,split('/')[3], index, sub_name, sub_desc,strip(), sub_link])
         39
                         index = index + 1
         /@kkw119 | 0 | 아키세라믹 | 아키세라믹의 브런치입니다. | /@hkp0402
         /@kkw119 | 1 | Dan Lee | 커피 좀 좋아해서 내려 마시고 있고, 지금까지 엔지니어도 하고 개발도 하고 지금은 기획하면서 개발팀 맡고
```

## 웹 크롤링 : 구독자/ 관심작가 리스트 데이터 수집 결과

Α	В	С	D	E	F
to_link	index	from_name	from_desc	from_link	
/@dimension-value	0	스텔라	데이터분석을 공부합니다	/@ivoryoh	0609
/@dimension-value	1	castellan	castellan의 브런치입니다.	/@hancas	tellan
/@dimension-value	2	KyeongMi	KyeongMi의 브런치입니다.	/@kyeong	mimfu8
/@dimension-value	3	양인환	양인환의 브런치입니다.	/@reay121	12
/@dimension-value	4	박구라	배우는 법을 배워가는 중입니다.	/@gurapa	rk
/@dimension-value	5	김은아	김은아의 브런치입니다.	/@happyr	naa
/@dimension-value	6	joo	HE.StoryFolio	/@storyfo	io
/@dimension-value	7	허성오	데이터 사이언스 연구소	/@castlefi	ve
/@dimension-value	8	수연	수연의 브런치입니다.	/@xy3on	
/@dimension-value	9	이민규	대기업부터 스타트업까지, 무한경쟁에서 ROI를	/@milares	oltimi
/@dimension-value	10	정동우	우물 밖 개구리.세상을 빗대어 사람을 봅니다.	/@jucie50	0
/@dimension-value	11	마지	거북이와 요구르트를 좋아합니다.	/@thj0228	
/@dimension-value	12	매튜캉	스타트업, 여행, 커피, 재즈를 좋아하는 20대 직	/@minwo	okang
/@dimension-value	13	이규봉	이규봉의 브런치입니다.	/@lee9bo	ng
/@dimension-value	14	손정웅	손정웅의 브런치입니다.	/@eviljiyo	
/@dimension-value	15	유민서	유민서의 브런치입니다.	/@canadi	1029
/@dimension-value	16	heren	에디터 큐레이터 읽고 쓰고 맛보고 즐기고	/@brunch	kyif
/@dimension-value	17	이광춘	이광춘의 브런치입니다.	/@hdecal9	99
/@dimension-value	18	한재원	IT 스타트업 회사 '당근마켓'의 마케터 한재원인	/@onea	
/@dimension-value	19	Sunyoung Melissa Kim	Sunyoung Melissa Kim의 브런치입니다.	/@sunyou	ngmelissa
/@dimension-value	20	Data Influencer		/@datainf	luencer
/@dimension-value	21	mdoong	mdoong의 브런치입니다.	/@lampyk	
/@dimension-value	22	박현지	박현지의 브런치입니다.	/@e2julie	
/@dimension-value	23	리노		/@lsdphjy	
/@dimension-value	24	천영훈	천영훈의 브런치입니다.	/@youngo	ldangi
/@dimension-value	25	JH	JH의 브런치입니다.	/@jihyunh	aneoqw
/@dimension-value	26	여니맘	여니맘의 브런치입니다.	/@einie98	
/@dimension-value	27	yoomyoom	yoomyoom의 브런치입니다.	/@winty33	3
/@dimension-value		, 박결	학결의 브런치입니다.	/@gyul61	
/@dimension-value	29	버건디	버건디 공간	/@classicb	



#### 웹 크롤링 : 개인별 프로필

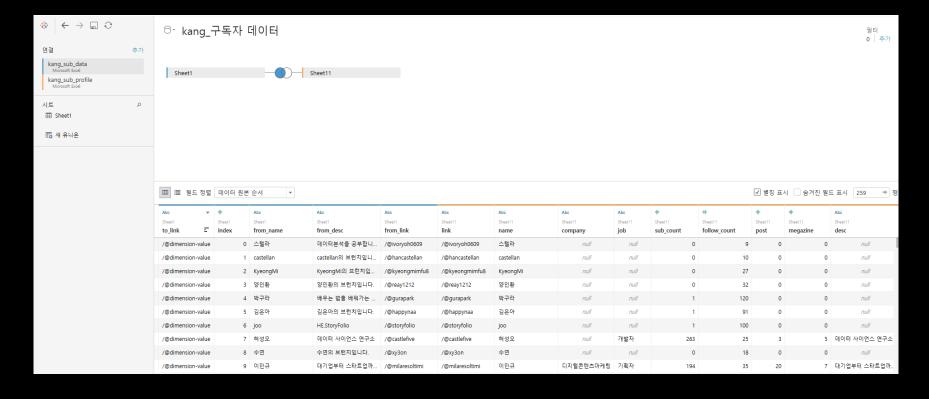


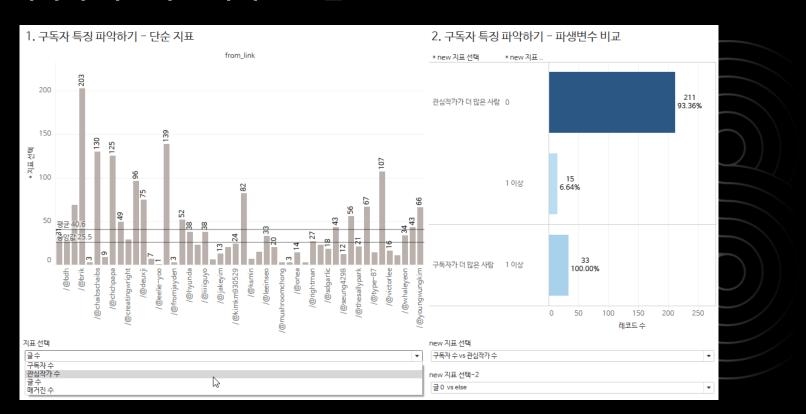


## 웹 크롤링 : 개인별 프로필 데이터 수집 결과

link	name	company	job	sub_count	llow_cour	post	megazine	desc	tag					
/@ivoryoh0609	스텔라			0	9	0	0							
/@hancastellan	castellan			0	10	0	0							
/@kyeongmimfu8	KyeongMi			0	27	0	0							
/@reay1212	양인환			0	32	0	0							
/@gurapark	박구라			1	120	0	0							
/@happynaa	김은아			1	91	0	0							
	joo			1	100	0	0							
/@castlefive	허성오		개발자	263	25	3	5	데이터 사이언스 연구소	IT,자기계팀	발,일,개발자	,연구자,강(	견자		
/@xy3on	수연			0	18	0	0							
	이민규	디지털콘텐츠미	기획자	194	35	20	7	대기업부터 스타트업까지, 무한경쟁	마케팅,콘	텐츠,브랜딩	,기획자,마	게터,강사,디	지털콘텐츠	드마케팅
/@jucie500	정동우			1	257	0								
/@thj0228	마지			0	6	0	0							
/@minwookang	매튜캉			0	24	0	0							
/@lee9bong	이규봉			0	6	0	0							
/@eviljiyo	손정웅			0	8	0	0							
/@canadi1029	유민서			0	23	0	0							
/@brunchkyif	heren			0	26	0	0							
/@hdecal99	이광춘			28	2851	0	0							
/@onea	한재원	당근마켓	마케터	31	12	14	3	IT 스타트업 회사 '당근마켓'의 마케	마케팅,스티	라트업,글쓰	기,마케터,	기획자,일러	스트레이터	,당근마켓
/@sunyoungmelissa	Sunyoung			0	54	0	0							
	Data Influ			0	24	0	0							
/@lampyk	mdoong			0	50	0	0							
/@e2julie	박현지			0	3	0	0							
/@lsdphjy	리노			0	38	0	0							
	천영훈			0	21	0	0							





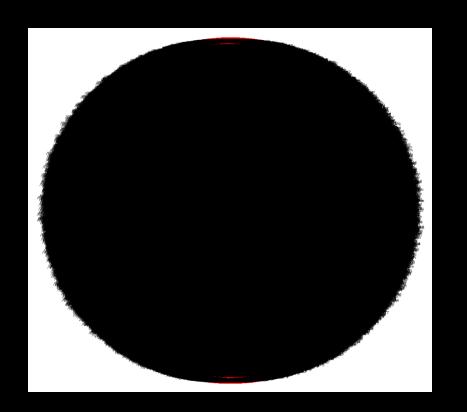


- 평균 구독자 수 90명
  - <u>구독자가 0명인 사람이 129명</u>, 이를 제외하면 평균 179명의 구독자 보유
- 평균 관심작가 수 135명
  - <u>관심작가가 1000명 이상인 사람</u>을 제외하면 평균 88명의 관심작가를 표시
- 평균 글 수 8개
  - <u>글을 1개 이상 발행한 사람은 48명 뿐</u>, 이들의 평균 글 수는 41개
- 평균 매거진 수 0개
  - 매거진이 0개인 사람을 제외하면 <u>평균 4개의 매거진 보유</u>

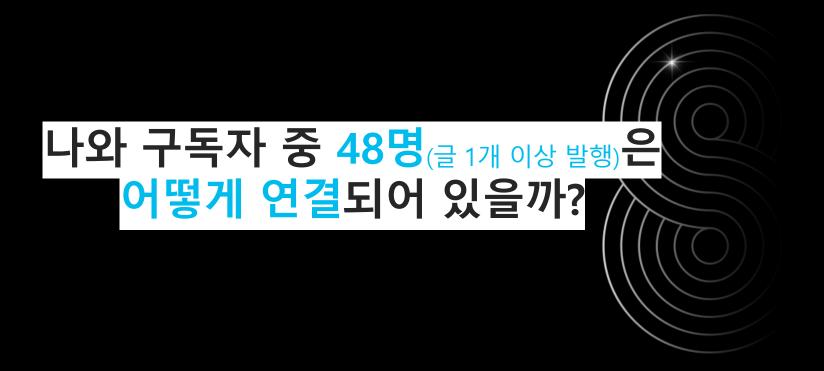




```
In [1]
        1 import pandas as pd
                                                                   import networkx as nx
                                                                   import matplotlib,pyplot as plt
In [7]
           data = pd,read_excel(',/data/kang_
                                                                   import sys
           data_head()
Out [7]
                                                   In [11]:
                                                                   E = nx,Graph()
             from link index
                              to name
       0 /@ivoryah0609
                              취준생LAB
                                                                   for word1, word2 in zip(from_list, to_list):
                                                                        E,add_edge(word1, word2)
       1 /@ivoryoh0609
       2 /@ivoryoh0609
                       2 Kim Koala Eddy El L
                                                   In [12]:
                                                                   # 항글 표시되도를 푼트 조정
       3 /@ivoryoh0609
                                빚그림
                                                                   if sys.platform in ["win32", "win64"]:
       4 /@ivoryoh0609
                                                                        font_name = "malgun gothic"
                                                                   elif sys.platform == "darwin":
                                                                        fornt_name = "AppleGothic"
       1 from_list = data['from_link']
          to_list = data['to_link']
                                                                   plt_figure(figsize=(20,20))
        1 for word1, word2 in zip(from_list
                                                                   plt.axis("off")
              print (word1.word2)
       /@ivoryoh0609 /@jobseekerlab
                                                                   G = nx, Graph(E)
                                                                                        # 일반적인 그래프, 모든 선에 대해 선을 이어줌
       /@ivorvoh0609 /@dimension-value
                                                                   # 8 = nx.minimum spannino tree(E) # Minimum spannino tree 선 군는 것 회소회 해서 전체를 이어쯤
       /@ivorvoh0609 /@hero4earth
       /@ivorvoh0609 /@bdh
       /@ivoryoh0609 /@mapthecity
                                                                   nx,draw_circular(G,
       /@ivoryoh0609 /@gimmesilver
                                                              14
                                                                            font_family=font_name,
       /@ivorvoh0609 /@cloud09
       /@ivorvoh0609 /@bonfire
                                                                             node_color="red".
       /@ivoryoh0609 /@brunch
                                                              16
                                                                           label_pos=0, #0=head, 0.5=center, 1=tail
       /@hancastellan /@elicecoding
                                                                            with_labels=True,
       /@hancastellan /@dimension-value
                                                                           font_size = 12
       /Mhancastellan /Mdrvishin
       /@hancastellan /@ekdrum
                                                              19
       /@hancastellan /@leeeeesh
                                                                   # plf. show()
       /@hancastellan /@reinitiate
       /@hancastellan /@holidavmemories
       /@hancastellan /@ryhynk
                                                                   # 010131 即學 对容整다
       /@hancastellan /@c4u
                                                                   plt,savefig("샘플,png")
       /@hancastellan /@brunch
                                                              24
       /@kyeongmimfu8 /@dimension-value
```







```
5/895
                            /@soks90
                                          /(@pruncnkxco
                                                                   /@soks90/@bruncnkxco
             57896
                                                                 /@imagineer/@brunchkxco
                          /@imagineer
                                          /@brunchkxco
             57897
                           /@prorange
                                          /@brunchkxco
                                                                  /@prorange/@brunchkxco
             57898
                                                                   /@eastgo/@brunchkxco
                            /@eastgo
                                          /@brunchkxco
             57899
                            /@steven
                                          /@brunchkxco
                                                                    /@steven/@brunchkxco
             57900
                       /@kangsunseng
                                          /@brunchkxco
                                                              /@kangsunseng/@brunchkxco
             57901
                                                                 /@hoonyqqq/@brunchkxco
                          /@hoonyqqq
                                          /@brunchkxco
             57902
                            /@brunch
                                          /@brunchkxco
                                                                   /@brunch/@brunchkxco
            57903 rows × 3 columns
In [307]:
                 cd = (total_data['join'] == total_data['join'][10])
              2 total_data[cd]
Out [307] :
                                   from link
             10 /@dimension-value /@iucie500 /@dimension-value/@iucie500
In [295]:
                 total_data = total_data.reset_index(drop=True)
                 |total_data = total_data.drop_duplicates('join')
In [311]:
                 df = pd.DataFrame(total_data)
                 df.columns = ['to_link','from_link','join']
                 df.to_excel("total_network_data.xlsx", index=True)
```

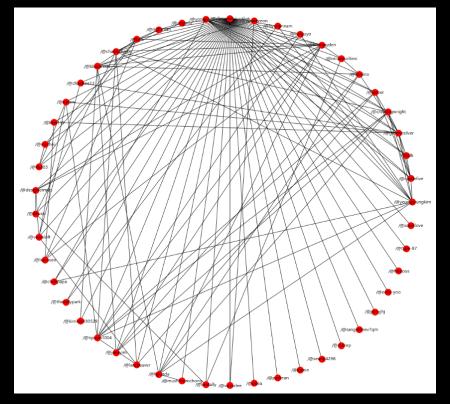
```
In [316]:
            1 import itertools
            3 case = list(map(''.join, itertools.permutations(id_list, 2)))
            4 case
             /@sangmileev/gm/@tipoon1UU4 .
            '/@sangmileev7qm/@hyunda',
            '/@sangmileev7qm/@sentally'
            '/@sangmileev7qm/@mushroomchong',
            '/@sangmileev7gm/@victorlee'.
            '/@sangmileev7gm/@thesallypark'.
            '/@sangmileev7qm/@latebeaver'.
            '/@sangmileev7qm/@whaleveon'.
            '/@sangmileev7gm/@visualoft'.
            '/@sangmileev7gm/@dimension-value'.
            '/@deuxii/@castlefive',
            '/@deuxii/@milaresoltimi'.
            '/@deuxii/@onea'.
            '/@deuxji/@youngwungkim',
            '/@deuxii/@bdh'.
            '/@deuxji/@parkean',
            '/@deuxii/@vongiiniinipIn'.
            '/@deuxji/@rightman'.
            '/@deuxji/@binjino',
            '/@deuxii/@gimmesilver'.
            /@do....; ; /@do2002;
```

```
In [320]: 1 cd1 = (total_data['join'] == case[2351])

Out[320]: to_link from_link join

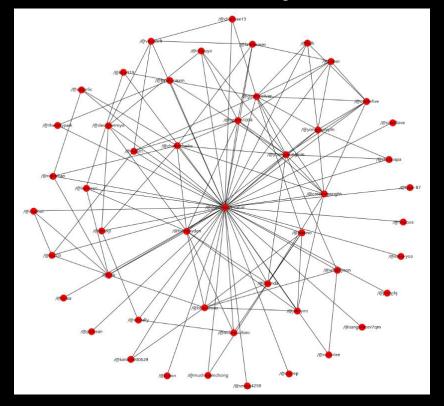
258 /@dimension-value /@visualoft /@dimension-value/@visualoft
```

```
In [375]:
             # 환글 표시되도록 푼트 조정
             if sys.platform in ["win32", "win64"]:
                 font_name = "malgun gothic"
             elif sys.platform == "darwin":
                 fornt_name = "AppleGothic"
             plt_figure(figsize=(20,20))
             plt_axis("off")
                               # 일반적인 그래프, 모든 선에 대해 선을 이어줌
             G = nx.Graph(E)
             #8 = nx.minimum_spanning_tree(E) # Minimum_spanning_tree 선 굿는 것 최소화 해서 전체를 이어움
          12
             nx.draw_circular(G.
         14
                     <del>font_family</del>font_name,
         15
                      node_color="red",
         16
                    label_pos=0, #O=head, O.S=center, 1=tail
                     with_labels=True,
         17
         18
                    font_size = 12
         19
             # plf. show()
```



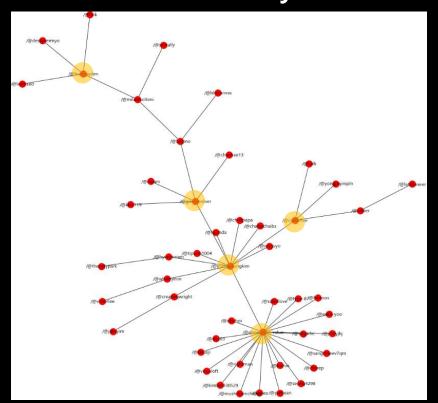


```
In [376]:
              # 환글 표시되도록 푼트 조경
             if sys.platform in ["win32", "win64"]:
                 font_name = "malgun gothic"
              elif svs.platform == "darwin":
                 fornt_name = "AppleGothic"
              plt_figure(figsize=(20,20))
              plt.axis("off")
             (G = nx,Graph(E)
                                  일반적인 그래프, 모든 선에 대해 선을 이어줌
              #U = nx.minimum_spanning_tree(E) # Winimum_spanning_tree 선 굿는 것 최소화 해서 전체를 이어움
              nx,draw_kamada_kawai(G,
                     <del>-font_family-fon</del>t_name,
          15
                      node_color="red",
          16
                     label_pos=0, #0=head, 0.5=center, 1=tail
          17
                     with_labels=True,
          18
                    font_size = 12
          19
              # plf. show()
```

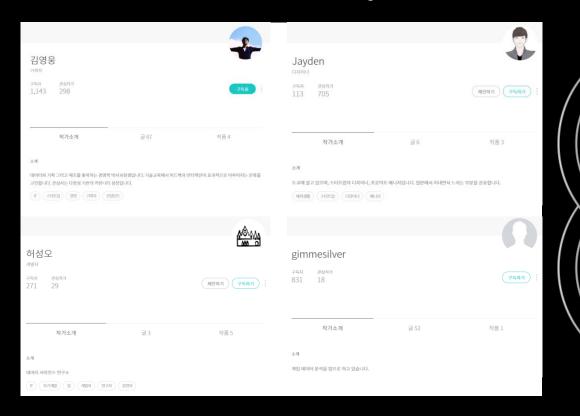




```
In [377]:
             # 항글 표시되도록 푼트 조경
             if sys.platform in ["win32", "win64"]:
                 font_name = "malgun gothic"
             elif sys,platform == "darwin":
                 fornt_name = "AppleGothic"
             plt_figure(figsize=(20,20))
             plt_axis("off")
             ## - mx Aceah(E) # 알반적인 그래프, 모든 선에 대해 선을 이어줌
             (G = nx,minimum_spanning_tree(E)
                                            #Minimum_spanning_tree 선 굿는 것 최소화 해서 전체를 이어줌
         12
             nx,draw_kamada_kawai(G,
                     font_family=font_name.
         15
                     node_color="red",
         16
                    label_pos=0, #O=head, O.5=center, 1=tail
         17
                     with_labels=True.
         18
                    font_size = 14)
         19
             # plf. show()
```



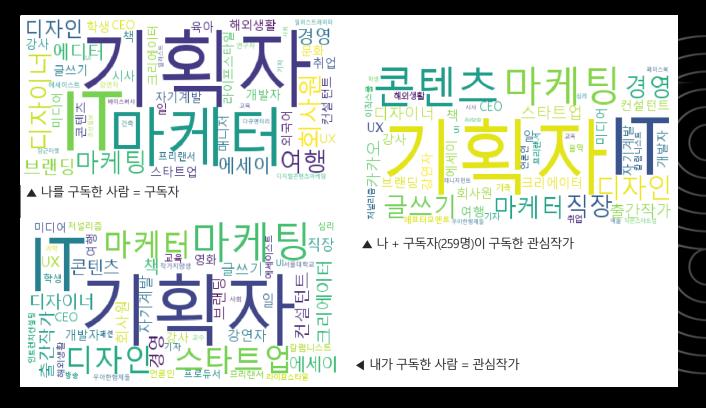




#### 데이터 시각화 기반의 인사이트 도출 with Python

# 우리는 비슷한 사람일까?





## 회사 밖의 데이터

04 프로젝트 후기



**SOSCON 2019** 

SAMSUNG OPEN SOURCE CONFERENCE 2019

#### 회사 밖의 데이터 프로젝트를 통해 얻은 결과

- 데이터를 활용해 궁금증을 직접 해결하기
- 현실적인 한계를 극복하기 위한 시도 -> 커뮤니티를 통한 학습과 성쟁
- 데이터 관련 컨퍼런스 등에서 콘텐츠로
   데이터 분석 과정에서 시각화의 역할 이야기하기



## THANK YOU



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019