

INF 558 Final Project:  
Korean-Pop  
WIKI

Xiaoyue Chen  
Qiuke Wang



# Objective

- Searching

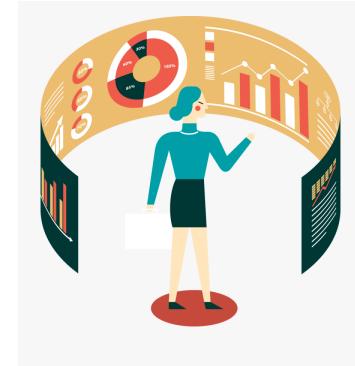
- Recommending



- Popularity Prediction



- Presenting statistic and trends with data visualization.



# Process & Tools



Data Crawling



BeautifulSoup



Web API

[import.io](#)



Mapping to  
ontology

[usc-isi-i2/rltk](#)



Entity Linking

Regular  
expression

String similarity



Knowledge  
Graph



Query & Web



1.0.0



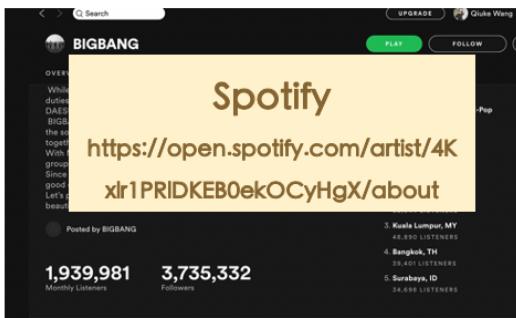
**Welsons Back** is the debut studio album by South Korean male group **ICON**. The full album was initially planned to be released on November 2, 2015. However, on October 27, 2015, **YG Entertainment** announced that the release of the full album would be delayed. This album was released on December 24, 2015, along with additional music via videos released on November 16.

## Discography

## Albums [edit]

Title	Album details	Peak positions	Sales
<a href="https://en.wikipedia.org/wiki/List_of_South_Korean_idol_groups_(2010s)">Wiki</a>			KOR: 204,249 <small>[27]</small>
"High High"		3	• KOR: 1,368
"Oh Yeah" (featuring Park Bom)	2010	2	• KOR: 1,38
"Knock Out"		5	• KOR: 775
"Zutte"	2015	2	• KOR: 1,01 • US: 6,000
			Album

Group label  
Group genres



## Monthly listeners

- Group name
- Group gender
- Group introduction
- Group members
- Group image
- Member information

# Data Sources

# Data Crawling

Challenge: inconsistent structure

```
><aside class="portable-infobox pi-background pi-theme-artist pi-layout-default">...</aside>
<b>BEBE6</b>
" (베베식스) is a six-member girl group under K12 Company. They debuted on August 10, 2017 with their first single "Shot Me"."
```

```
</aside>
<p>
<b>Big Brain</b>
" (빅브레인) is a four-member R&B ballad group under "
<a href="/wiki/MAJOR9" title="MAJOR9">MAJOR9</a>
". They have been active since 2010 and officially debuted on
October 22, 2015 with the single "Billionaire Sound" and their
breakthrough song "Welcome".
"
```

## Members

- Sang Hoon (상훈)
- Jin Yong (진용)
- Byeong Eun (병은)
- Hong Hyun (홍현)

## Members

Name	Position(s)	Year(s) active
Donghyun (동현)	Leader, Lead Vocalist	2011–2019
Hyunseong (현성)	Main Vocalist	2011–2019

## Solution:

Optimizing spider and take more cases into consideration.

# Data Crawling

Challenge: missing information

Group gender: Girl group or Boy group or Co\_ed group? Number of group member?  
How to measure popularity?

## Solution:

Extracting desired information from unstructured data(introduction) or inferring from known knowledge.

BLACKPINK (블랙핑크; stylized as BLACKPINK) is a four-member girl group under YG Entertainment. They debuted on August 8, 2016 with their digital single album "Square One".

$$\text{Count}(\text{members}) = 4$$

Gender = 'Female'

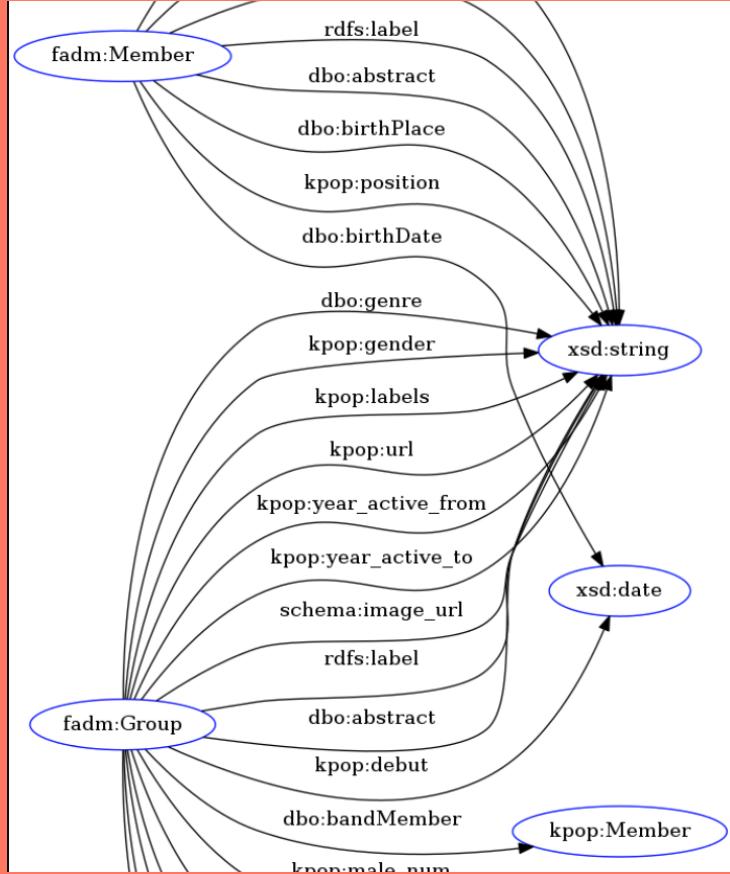
Accuracy: 0.880000

Precision: 0.916000

Recall: 0.880000

F1 score: 0.894975

Name	Position(s)	Years active
Jisoo (지수)	Lead Vocalist, Visual	2016–present
Jennie (제니)	Main Rapper, Vocalist	2016–present
Rosé (로제)	Main Vocalist, Lead Dancer	2016–present
Lisa (리사)	Main Dancer, Lead Rapper, Sub Vocalist, Maknae	2016–present



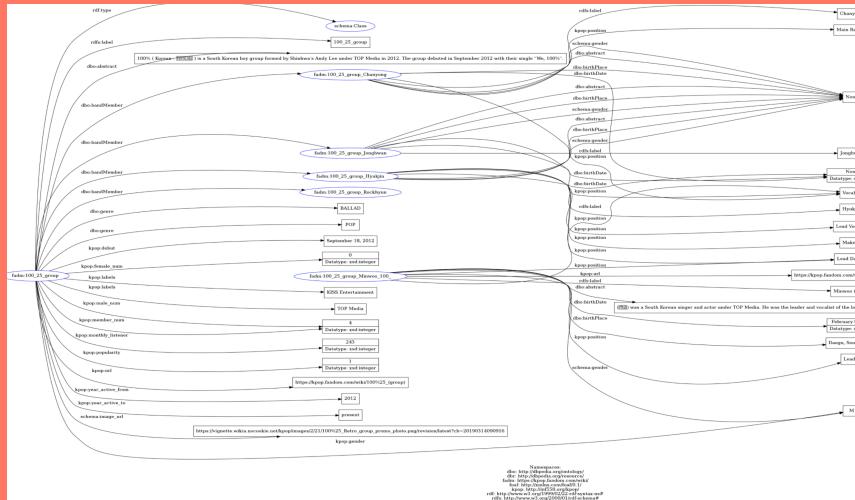
> image

```

1 @prefix dbo: <http://dbpedia.org/ontology/> .
2 @prefix dbr: <http://dbpedia.org/resource/> .
3 @prefix fadm: <https://kpop.fandom.com/wiki/> .
4 @prefix foaf: <http://xmlns.com/foaf/0.1/> .
5 @prefix kpop: <http://inf558.org/kpop/> .
6 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
7 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
8 @prefix schema: <http://schema.org/> .
9 @prefix xml: <http://www.w3.org/XML/1998/namespace> .
10 @prefix xsd: <http://www.w3.org/2001/XMLSchema#> .

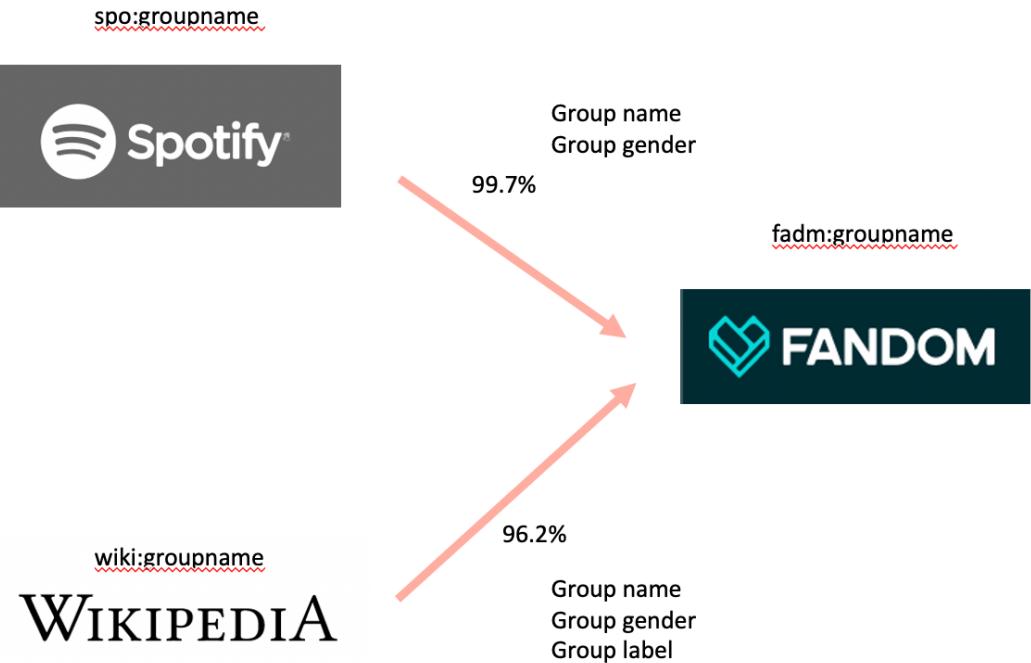
11
12 <https://kpop.fandom.com/wiki/100\_25\_group> a schema:Class ;
13   rdfs:label "100_25_group" ;
14   dbo:abstract "100% ( Korean : 백퍼센트 ) is a South Korean boy group formed by Shinhwa's And
15   dbo:bandMember <https://kpop.fandom.com/wiki/100\_25\_group\_Chanyong>,
16     <https://kpop.fandom.com/wiki/100\_25\_group\_Hyukjin>,
17     <https://kpop.fandom.com/wiki/100\_25\_group\_Jonghwan>,
18     <https://kpop.fandom.com/wiki/100\_25\_group\_Rockhyun> ;
19   dbo:genre "BALLAD",
20     "POP" ;
21   kpop:debut "September 18, 2012" ;
22   kpop:female_num 0 ;
23   kpop:gender "M" ;
24   kpop:labels "KISS Entertainment",
25     "TOP Media" ;
26   kpop:male_num 4 ;
27   kpop:member_num 4 ;
28   kpop:monthly_listener 243 ;
29   kpop:popularity 1 ;
30   kpop:url "https://kpop.fandom.com/wiki/100%25\_\(group\)" ;
31   kpop:year_active_from "2012" ;
32   kpop:year_active_to "present" ;

```



# Ontology

# Entity Linking



Kpop:company ----- rdfs:label

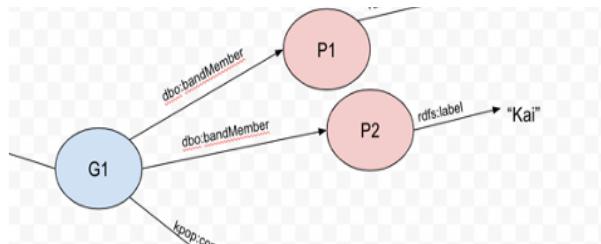
BNM

Brand New Music

브랜뉴 뮤직

BrandNew Music

Using dictionary to direct convert the different name to the same label.



rdfs:label ----- fandm:groupname

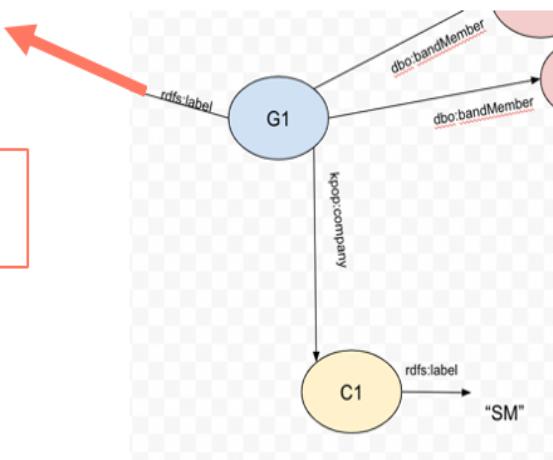
G-idle

(G)I-DLE

GI-DLE

GIDLE

GIDLE(group)



Lower case, delete all unnecessary punctuation, using gender label.

Levenshtein similarity > 0.9

# Entity Resolution

# SPARQL

## SPARQL query

```
Please input your sparql query here  
SELECT distinct ?group ?groupname ?imageurl  
WHERE{  
?group a schema:Class.  
?group schema:image_url ?imageurl.  
?group rdfs:label ?groupname.  
?group kpop:gender ?gender.  
filter regex(?gender,'M')  
?group kpop:labels ?company.  
filter regex(?company,'')  
?group kpop:member_num ?number.  
filter (?number)  
}  
LIMIT 5
```

QUERY

groupname	imageurl
TVXQ!	https://vignette.wikia.nocookie.net/kpop/images/d/d4/TVXQ%21_XV_promo_photo_1.png/revision/latest?cb=20191104063447
SuperM	https://vignette.wikia.nocookie.net/kpop/images/2/2c/SuperM_Let%27s_Go_Everywhere_promo_photo_1.png/revision/latest?cb=20191119001756
SHINee	https://vignette.wikia.nocookie.net/kpop/images/d/d8/SHINee_Sunny_Side_promo_photo.png/revision/latest?cb=20180706154531
SUPER_JUNIOR	https://vignette.wikia.nocookie.net/kpop/images/f/fa/SUPER_JUNIOR_Timeless_group_promo_photo.png/revision/latest?cb=20200115010831
EXO	https://vignette.wikia.nocookie.net/kpop/images/l/16/EXO_Obsession_group_concept_teaser_photo_1.png/revision/latest?cb=20191127040511

```
if num == '':
    queryline = "SELECT distinct ?group ?name WHERE{ ?group a schema:Class. ?group kpop:labels ?company.filter regex(?com
else:
    queryline = "SELECT distinct ?group ?name WHERE{ ?group a schema:Class. ?group kpop:labels ?company.filter regex(?com
sparql.setQuery(prefix + queryline)
temp = sparql.query().convert()
resultGroup = []
if len(temp["results"]["bindings"]) > 0:
    keysGroup = temp["results"]["bindings"][0].keys()
    for i in range(len(temp["results"]["bindings"])):
        line = []
        for key in keysGroup:
            if temp["results"]["bindings"][i][key]["type"] == 'uri':
                #need to replace link
                line.append([temp["results"]["bindings"][i][key]["value"],True])
            else:
                line.append([temp["results"]["bindings"][i][key]["value"],False])
        resultGroup.append(line)
```

# Choosing KG Platforms

- Apache Jena provides Fuseki components.
- Fuseki is a SPARQL server that supports the SPARQL language for retrieval and runs efficiently on both stand-alone and server sides.
- For the web application we built, we do need more query functions than the visualization.

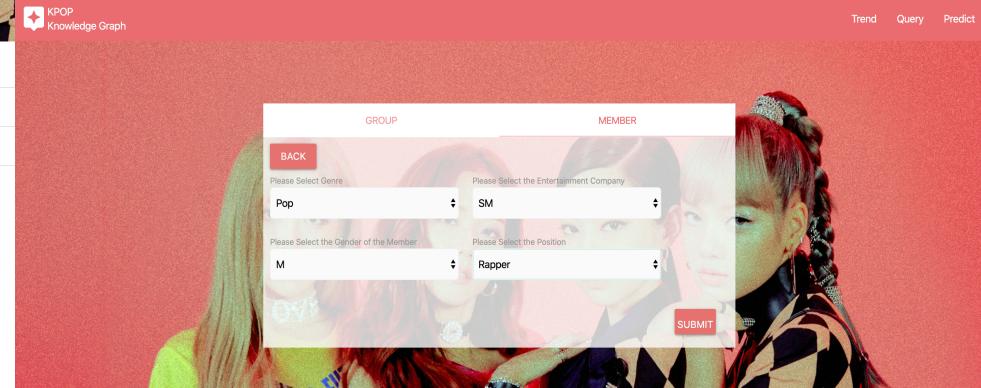
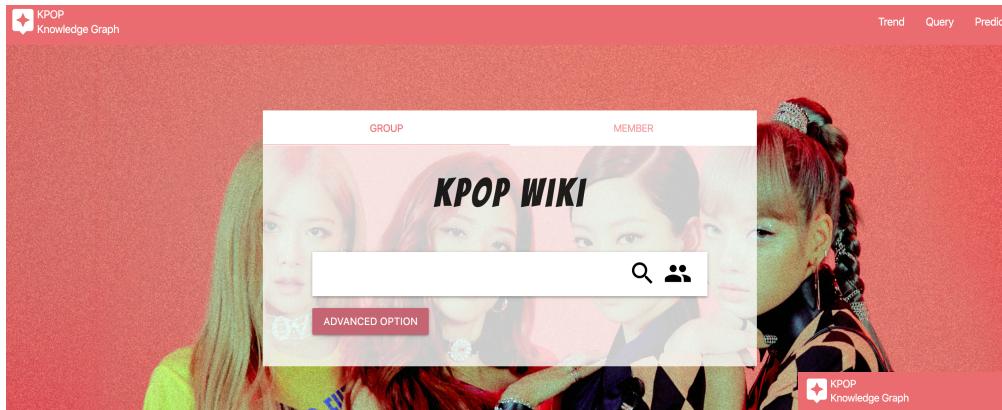


VS

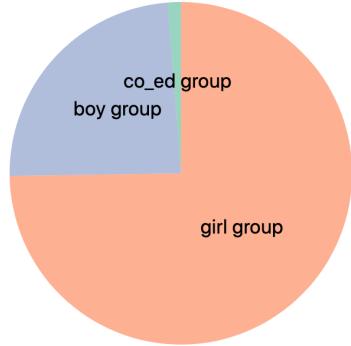


# Web

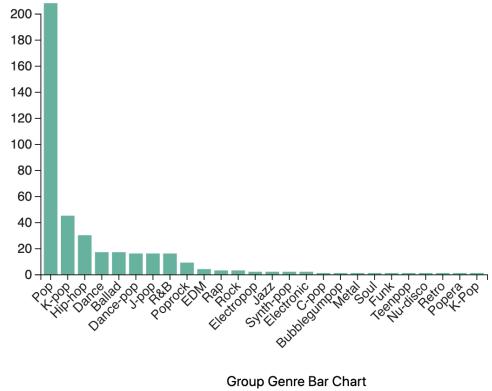
# search & detail Page



# Trend Page

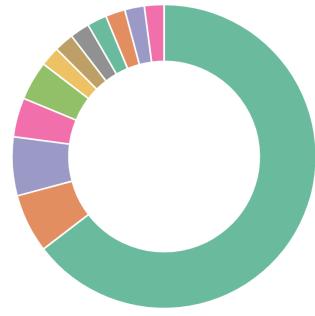


Group Gender Pie Chart

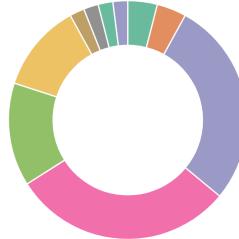


Group Genre Bar Chart

Group  
Overview



top 50 Most Popular K-Pop Groi



top 50 Most Popular K-Pop Groups Number of Members

Trend  
Overview

**TOP 10 MOST POPULAR K-POP GROUPS (2020)**

	<b>1</b> BLACKPINK		<b>2</b> BTS
	<b>3</b> EXO		<b>4</b> Wanna_One
	<b>5</b> GOT7		<b>6</b> SuperM
	<b>7</b> BIGBANG		<b>8</b> SEVENTEEN
	<b>9</b> Girls_27_Generation		<b>10</b> TWICE

# Predict Popularity

Select Gender

Girl Group  Boy Group  Co\_ed Group

Select Genres(multiple)

Choose group genres

Select Company

Choose company

Select number of members

Choose number of members

**PREDICT**

**Predicted Popularity:**

★★★★★

**Similar Group:**

 The\_Grace

 Girls\_27\_Generation

 S.E.S.

How to predict popularity when we can't find similar group?

## Solution:

Embedding the features(genre, gender, company and number of members). Train SVM + linear regression model. Query the avg popularity basing on different features, then add the predict score of our model. Using voting method to get the final score.

Machine  
Learning

+

Find  
Similarity

+

SPARQL

# Evaluation



Data Crawling



Mapping to ontology



Entity Linking



Query & Web

Reviewing the data scraping from the web. Analyzing the wrong data and improve the spider.

Using RDF validation tool.

Labeling the test data manually and detect the accuracy.

Using SPARQL to test on the data and validate manually.

# Improvement

We invited friends to help us test the whole application. We built a more user-friendly project according to their feedback.

GROUP: <https://kpop.fandom.com/wiki/TXT>.  
BIRTHPLACE: Honolulu, Hawaii .  
POSITION: Center , Maknae , Visual , Lead Vocalist .  
URL: <https://kpop.fandom.com/wiki/Hueningkai> .

## 1. Add more hyper-links

In the member detail page, we add not only the attributes belonging to the member, but also the group information the member belongs to.



## 3. Retrieve more information

In order to present the images of each member and group., we crawled the data source again to make the description page looks better.

```
https://your_group_member_detail_page  
schema:image_url "https://vignette.wikia.nocookie.net/  
ENTERTAINMENT/.../image/.../image.jpg"
```

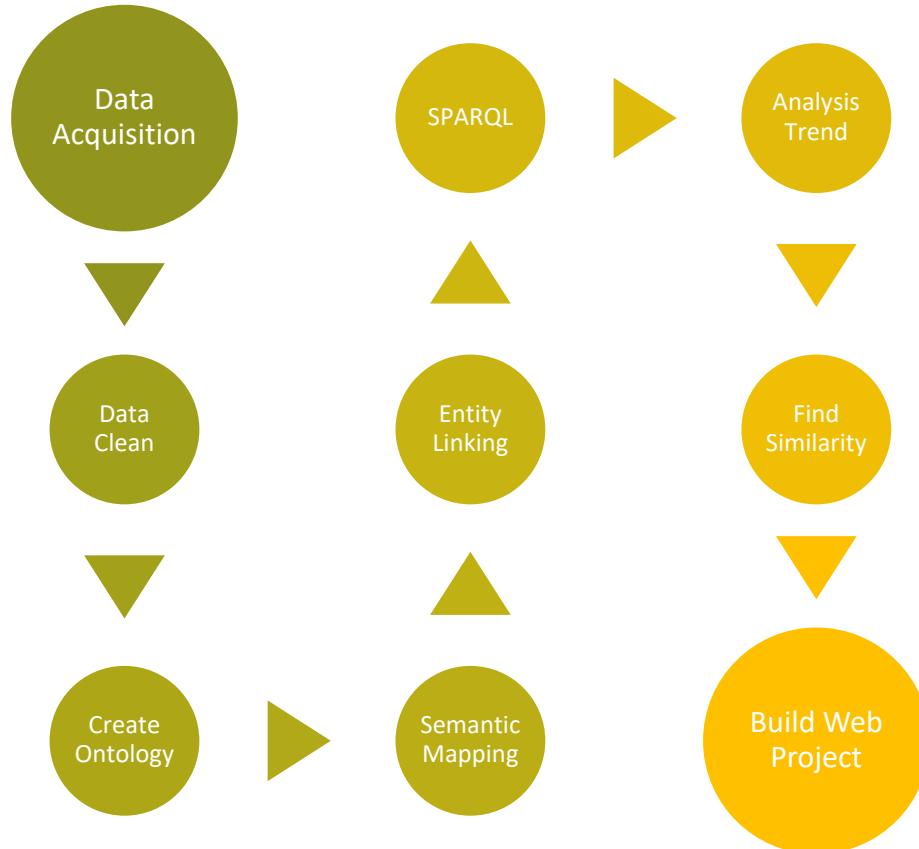
MAMAMOO



## 2. Change UI

After listening to users' advice, we combined the group search and member search together into one page using tab.

# Conclusion



The system helps us to digest all the knowledge we learnt from the class. From crawling data to building a full-stack web project.



We learnt that knowledge graph is a better way when trying to analyze complicate relationships among different entities.



With the trend page and predict tool, we have a better understanding of data visualization and machine learning.

A photograph of seven young men of diverse ethnicities, all looking directly at the camera from a low angle. They are wearing casual clothing, including t-shirts and hoodies, in various colors like white, blue, and grey. The background is a plain, light color.

Thank  
You for  
Watching!