



Similarity in Wikipedia Articles

Badenes, Carlos (*cbadenes*)

Garijo, Daniel (*dgarijo*)

Priyatna, Freddy (*fpriyatna*)

{*}@fi.upm.es

EDBT Summer School 2015

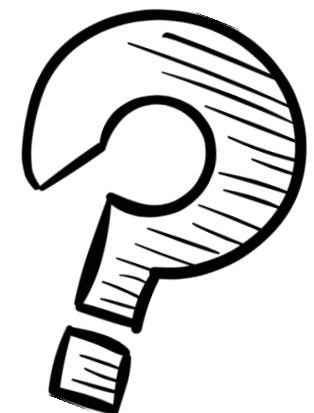


The screenshot shows the Wikipedia article for Madrid. At the top, there's a navigation bar with 'Article' and 'Talk' tabs, and a search bar. Below this is a banner for 'Wiki Loves Monuments'. The main heading is 'Madrid', followed by the text 'From Wikipedia, the free encyclopedia'. A coordinate box shows 'Coordinates: 40°23'N 3°43'W'. The article text begins with a disambiguation note: 'This article is about the capital city of Spain. For the autonomous community, see *Community of Madrid*. For other uses, see *Madrid (disambiguation)*.' The main text describes Madrid as a south-western European city and the capital and largest municipality of Spain, with a population of almost 3.2 million. It mentions its location on the Manzanares River and its status as the political, economic, and cultural center of Spain. A section on the Madrid urban agglomeration follows, noting its third-largest GDP in the European Union. The article also mentions its role as a major financial center and its status as one of the world's most livable cities. A sidebar on the left contains various navigation links like 'Main page', 'Contents', 'Featured content', and 'Tools'. At the bottom of the article, there are images of the Madrid flag and coat of arms.

Wikipedia Article:

- text
- links
- categories

Similarity between Wikipedia Articles





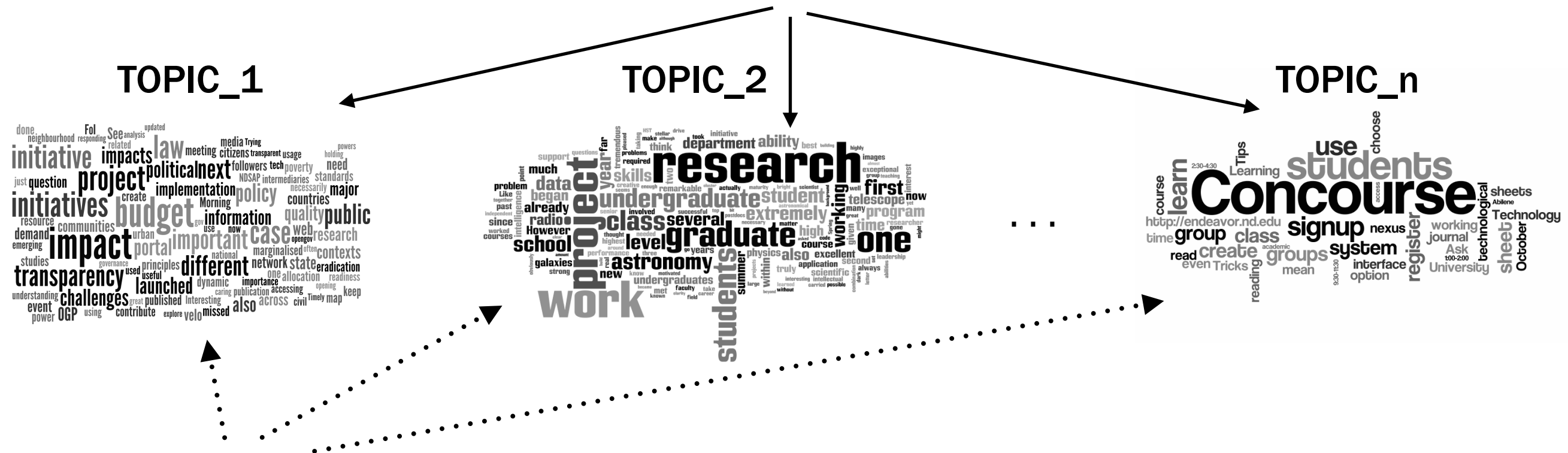
The screenshot shows the Wikipedia article for Madrid. At the top, there's a navigation bar with 'Article' and 'Talk' tabs, and a search bar. Below this is a banner for 'Wiki Loves Monuments'. The article title 'Madrid' is prominently displayed, followed by a coordinate string 'Coordinates: 40°23'N 3°43'W'. The main text begins with a disambiguation note and then describes Madrid as the capital of Spain, mentioning its population and metropolitan area. To the right of the text is a gallery of images, including a panoramic view of the city, a street view, and the Royal Palace. Below the gallery are the Madrid flag and coat of arms. The left sidebar contains various Wikipedia navigation links like 'Main page', 'Contents', and 'Tools'.

Wikipedia Article:

- text $\alpha \cdot simText$
- +
- links $\beta \cdot simLinks$
- +
- categories $\gamma \cdot simCtg$

$$simWA(R1, R2) = \alpha \cdot simTtxt(R1, R2) + \beta \cdot simLinks(R1, R2) + \gamma \cdot simCtg(R1, R2)$$

where $\alpha + \beta + \gamma = 1$



No matter the scale or the speed of development show that domestic private bank is still at the primary stage. Domestic commercial bank developed intermediate business again since 2000 from the starting point of view. China Merchants Bank launched Golden Sunflower financing as the symbol in 2003, and then, other banks launched their new financial product one after another. Personal finance, VIP wealth management, etc. have been developed rapidly. Domestic commercial banks have accumulated a large number of customers, and have tried to separate customer groups accurately, and thus consciously foreshadow the further development of private banking business.

Ri

HSBC is one of the world's largest banking and financial services organizations. Since opening in Shanghai in 1865, HSBC's business has grown rapidly in China, especially as the trade financing. In 1997, HSBC is one of the largest foreign banks in China, with 100 branches in Pudong district in Shanghai. In 1998, HSBC became the first-class member of the national RMB interbank lending market, and was approved to buyback and trade RMB through the market. HSBC China opened on April 2, 2007, head office in Shanghai. HSBC China has a total of 100 foreign bank expand and develop branches the most rapidly by mode of monopolization or investment that set foot in some areas including trade, bank, security and insurance. HSBC China has 100 branches in Shanghai Pudong New Area village bank in China. In addition, HSBC China is one of the largest investment foreign banks in the mainland, the total of investment is more than US\$5 billion, including Bank of Shanghai, HSBC, Citibank, etc.

From private banks at the top of the pyramid to the bottom of the rural finance, from single bank business to a necessary condition for mixed operation such as insurance, trust and securities, HSBC is always the platform and partner of universal bank have been formed.

Rj

$$sim_{cont}(R_i, R_j) = 10^{-JSD(p,q)}$$

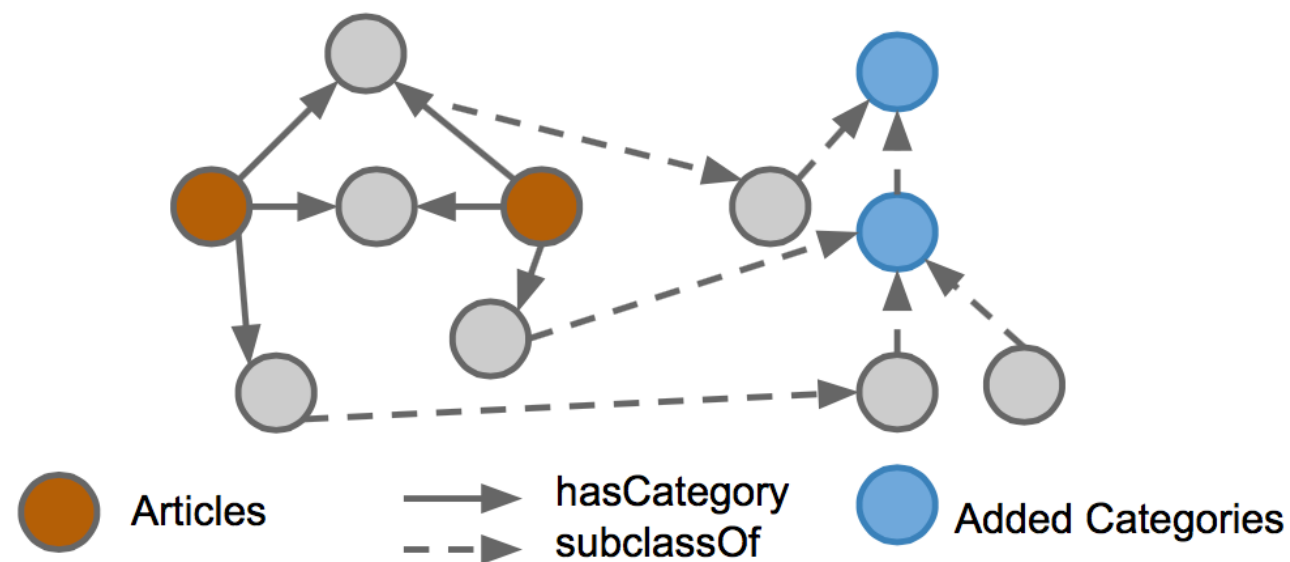
$\mathbf{p} = [0.5, 0.3, \dots, 0.7]$

$$q = [0.2, 0.4, \dots, 0.9]$$



Articles with multiple common **categories** are likely to be similar

$$\text{Sim}(A,B) = \text{cat}(A) \cap \text{cat}(B) / ((\text{cat}(A) \cup \text{cat}(B)) / 2)$$



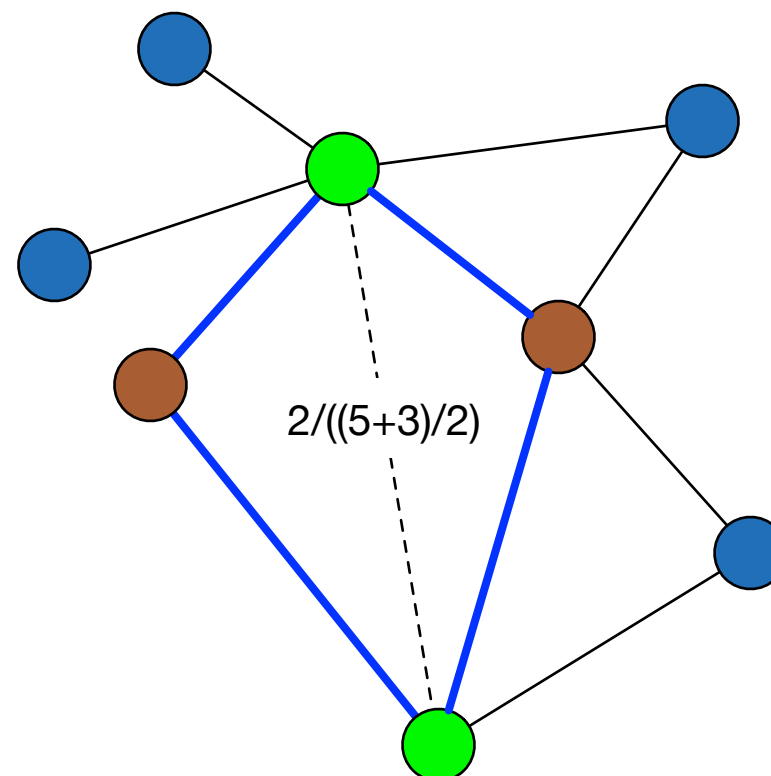
Noise filtering is necessary (e.g., “All articles lacking in-text citations”).

See https://github.com/cbadenes/siminwikart-challenge4/blob/master/category/wikipedia_bad_categories.txt



Articles with multiple common **links**
are likely to be similar

$$\text{Sim}(A,B) = \text{links}(A) \cap \text{links}(B) / ((\text{links}(A) \cup \text{links}(B)) / 2)$$



(simLinks) $\alpha = 0.2$
 (simCtg) $\beta = 0.2$
 (simTxt) $\gamma = 0.6$

Fernando Alonso



simTxt = 0.065
 simLinks = 0.0
 simCtg = [1]0.095
 [3]0.161

simTxt = 0.059
 simLinks = 0.019
 simCtg = [1]0.117
 [3]0.181

[1]0.062
 [3]0.075

[1]0.068
 [3]0.069

simTxt = 0.052
 simLinks = 0.019
 simCtg = [1]0.166
 [3]0.172

[1]0.019
 [3]0.023

Lionel Messi



[1]0.666
 [3]0.683

[1]0.043
 [3]0.072

simTxt = 0.980
 simLinks = 0.175
 simCtg = [1]0.217
 [3]0.302

simTxt = 0.060
 simLinks = 0.008
 simCtg = [1]0.030
 [3]0.172



Iker Casillas

simTxt = 0.069
 simLinks = 0.004
 simCtg = [1]0.080
 [3]0.134

[1]0.058
 [3]0.069



Princess Akiko

Lionel Messi

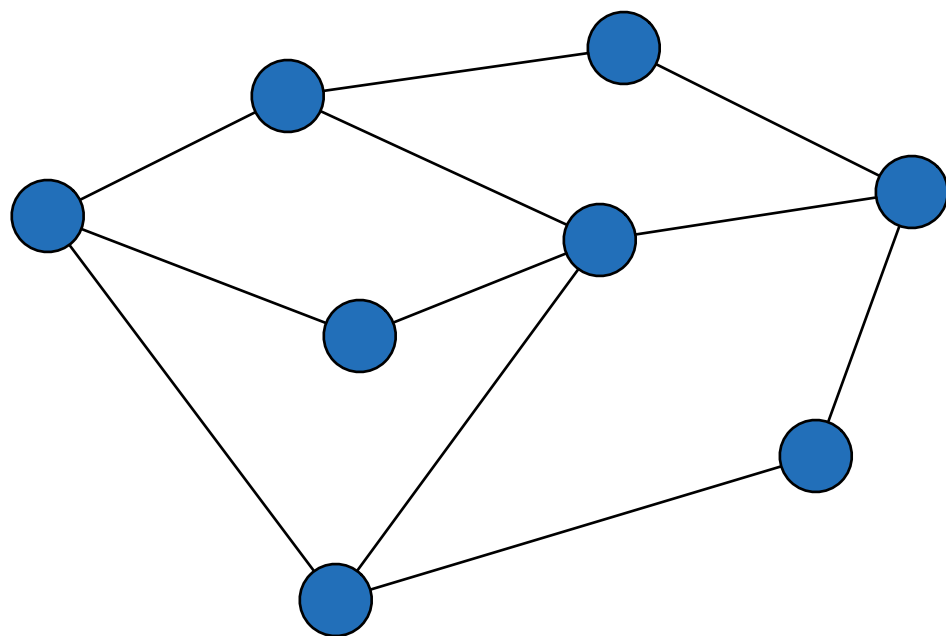


simTxt = 0.060 -> <common words>
simLinks = 0.008 -> (England, Buenos_Aires, Chile, Madrid, Argentina)
simCtg=[1]0.030 -> living_person

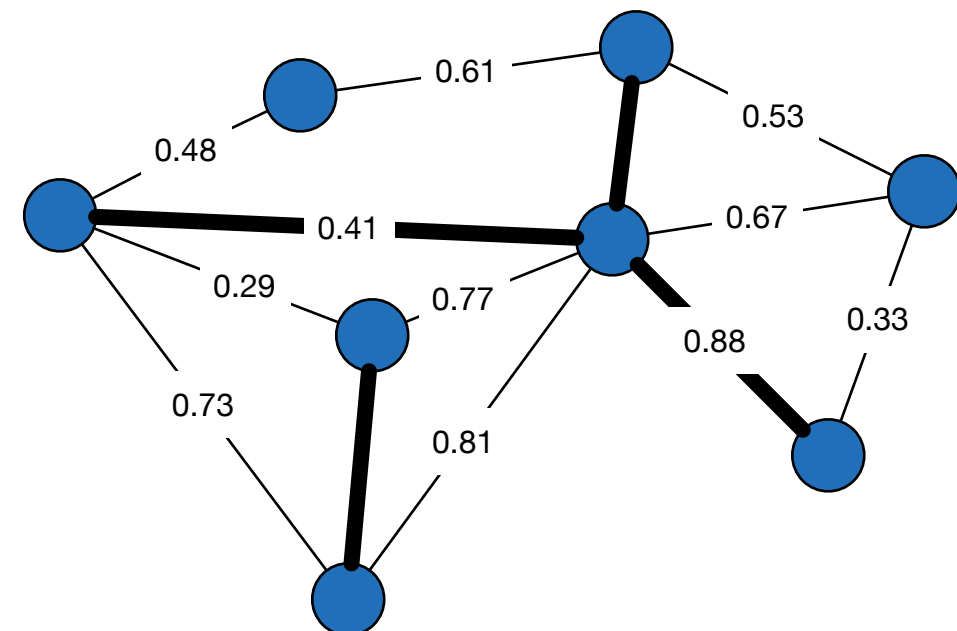


Princess Akiko

Graph based on Links



Graph based on Similarities

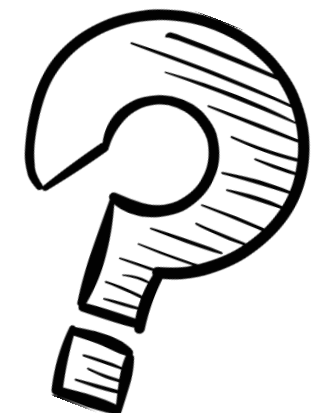




The screenshot shows the Wikipedia article for Madrid. At the top, there's a navigation bar with 'Article' and 'Talk' tabs, and a search bar. Below this is a banner for 'Wiki Loves Monuments'. The article title 'Madrid' is prominently displayed, followed by a coordinate line. The main text begins with a disambiguation note and then describes Madrid as the capital of Spain, mentioning its population and metropolitan area. To the right of the text is a gallery of images, including a cityscape, a street view, and the Royal Palace. Below the gallery are the Madrid flag and coat of arms. The left sidebar contains various Wikipedia navigation links like 'Main page', 'Contents', 'Featured content', and 'Tools'.

Wikipedia Article:

- text
- links
- categories



Wikipedia links reliability
(missing links)

Similarities between categories (as topics)
can define relations between articles



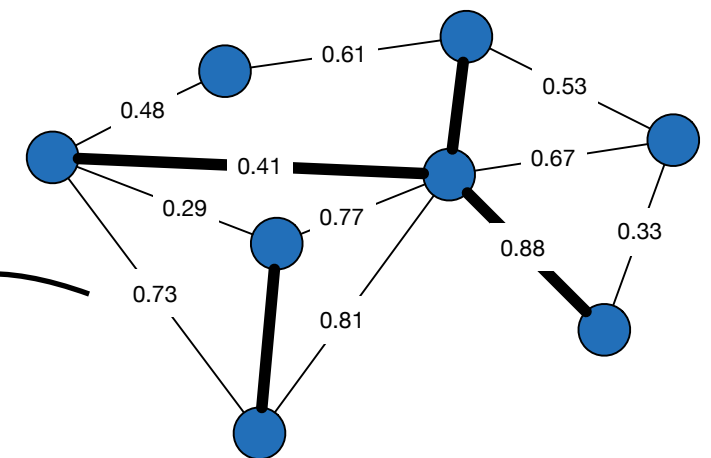
Graph based on Links

Latent
Dirichlet
Allocation

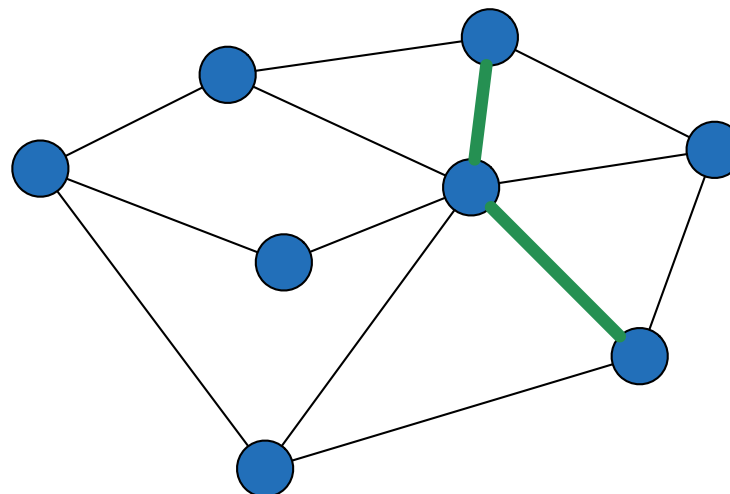
Topic Model



Graph based on Similarities



Subgraph Pattern Matching



GitHub



<https://github.com/cbadenes/siminwikart-challenge4>

Spark



Scala

*Sparksee