

Topic-Based Annotations

Up to now we always used a term-based vector space model, where we use the cosine between a query and a document to measure the similarity. However this approach has some issues if we consider synonymy and polysemy in a document or in a query.

EXAMPLE

"MICROSOFT'S BROWSER"

SYNONYMY

"INTERNET EXPLORER"

DIFFERENT TERMS, SAME MEANINGS

"PAPARAZZI PHOTOMORPHING STARS"

POLYSEMY

"ASTRONOMER PHOTOGRAPHING STARS"

SIMILAR PHRASES, DIFFERENT MEANINGS

We would like to be able to retrieve

the information needed by the user.

We do this using a knowledge graph,

an hand-crafted graph where the nodes represent entities and the edges represent relations.

A knowledge graph contains context-aware informations

and the idea of the topic-based annotation

is to find the anchors in the text

and associate them with the correct entities

drawn from this graph.

THE MOST KNOWN KNOWLEDGE GRAPH
IS WIKIPEDIA

EXAMPLE

"MICROSOFT'S BROWSER"



| E
WIKIPEDIA
PAGE

"INTERNET EXPLORER"



| E
WIKIPEDIA
PAGE

"PAPARAZZI PHOTOGRAPHING STARS"



| E
WIKIPEDIA
PAGE

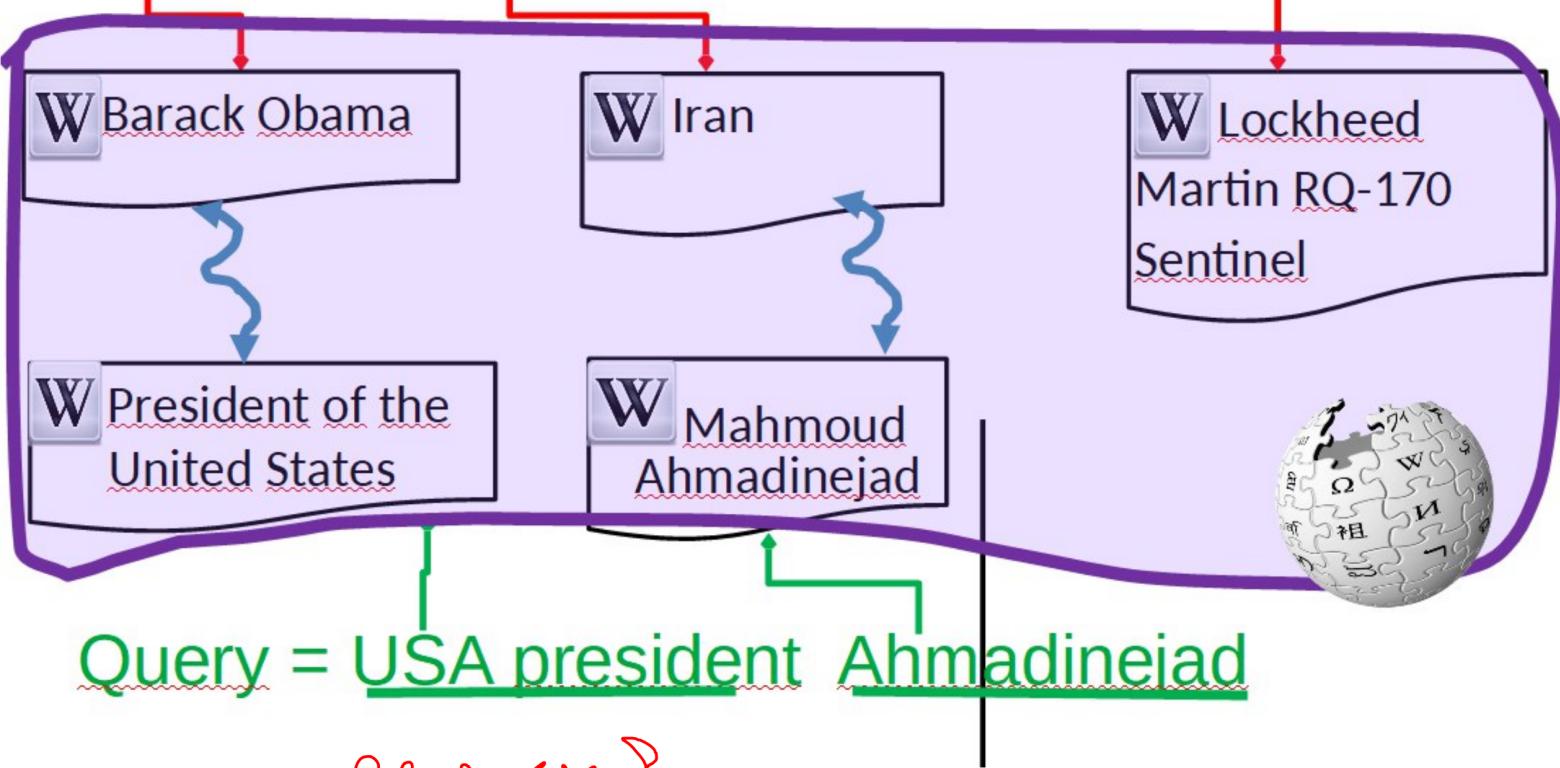
"ASTRONOMER PHOTOGRAPHING STARS"



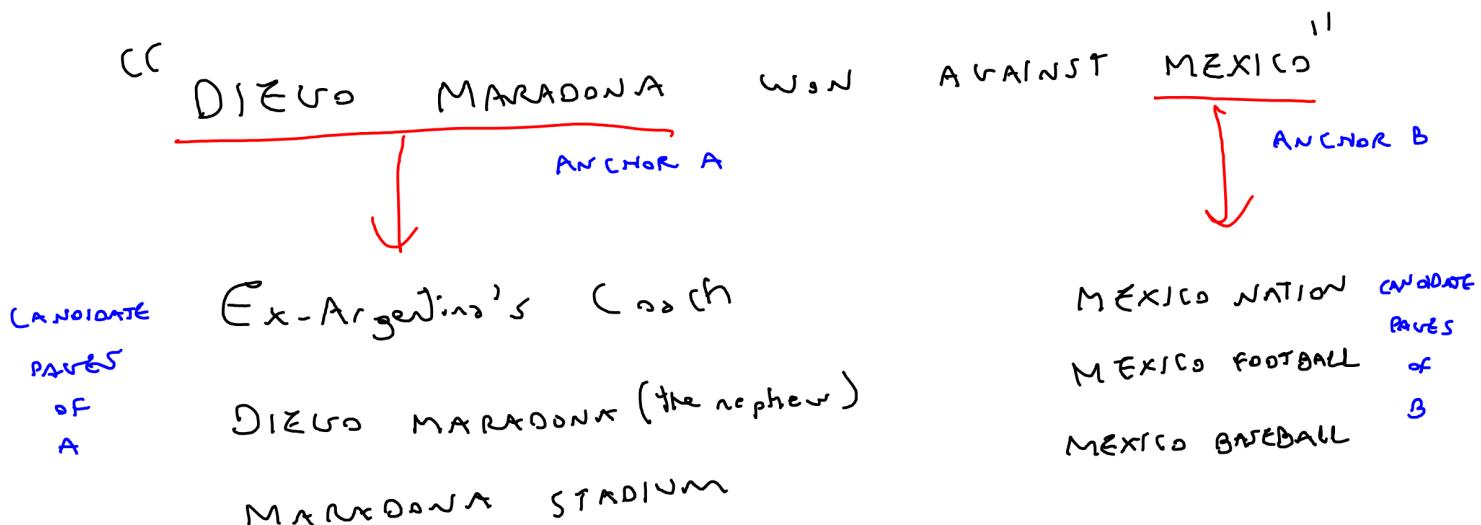
| E
WIKIPEDIA
PAGE

EXAMPLE OF QUERY

"Obama asks Iran for RQ-170 sentinel drone back"



BUT PROBLEM ?



How do we pick the correct entities (wiki pages) for our two anchors? Clearly they are related and the key is to quantify the relatedness between their candidate pages.

So how do we map our two anchors to the correct Wikipedia pages?

NOTE we are assuming that two anchors are enough to disambiguate the context (for simplicity)

Here what we need to link an anchor a with the correct page p :

- Context of a

Given the text T , the context of a are all the words around the anchor

$$T = \dots w_1 w_2 w_3 \dots w_4 w_5 w_6 \dots$$

- Context of p

$$\text{Simply the context } p = z_1 z_2 z_3 \dots$$

- Link probability measure

Measures the "goodness" of an anchor a (in how many Wikipedia pages it appears as a link?)

$$L_p(a) = \frac{\text{freq. of } a \text{ as anchor}}{\text{freq. of } a \text{ in text}}$$

- Commonness measure

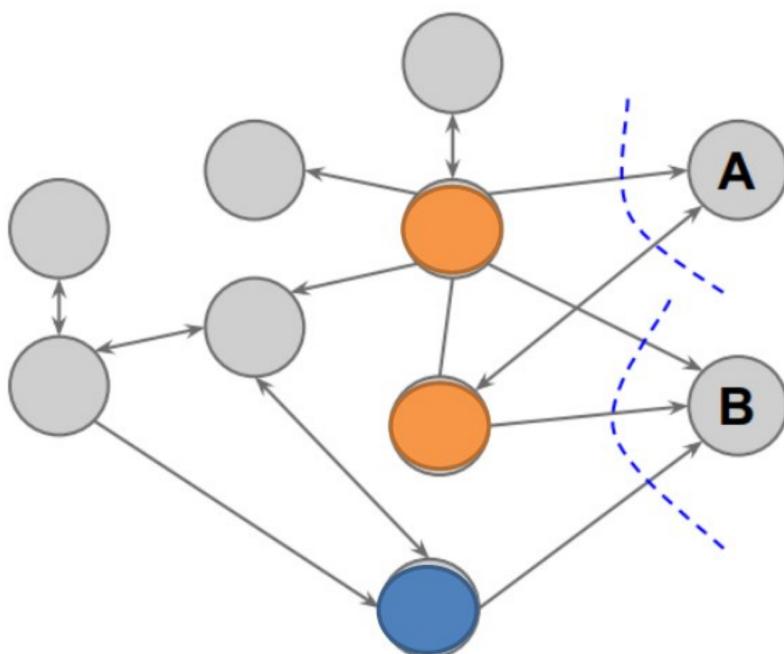
Measures how many times a certain anchor a links to the same page P

$$P(p|a) = \frac{\# \text{a linked to } p}{\# \text{a as anchor}}$$

- Relatedness between pages

In the knowledge graph, p is a entity node and a is a mention node

$$\text{sim}(A, B) = \frac{\log(\max(|A|, |B|)) - \log(|A \cap B|)}{\log(|\omega|) - \log(\min(|A|, |B|))}$$



Let's see now a concrete example of an annotator, **TagMe**, and see how it checks for possible anchors

Tag Me

A tool to identify meaningful short-phrases in a free text and link them to a Wikipedia page.

= It analyses the text, obtain a set of anchor and associate them with Wikipedia pages

By using the features described above, Tagme is able to disambiguate between the pages through a voting scheme

TagMe steps

1) Given \rightarrow text, parse it and
select the texts to annotate (anchors)

2) Retrieve the candidate pages for each
anchor

3) Dissambiguate & Prune pages:

- Prune by commonness $\leq T$
(T tradeoff: speed vs recall)

- Vote every candidate page, for both anchors

$$\text{Vote}_b(p_a) = \frac{\sum_{p_b \in P_g(b)} \text{sim}(p_a, p_b) \cdot P(p_b | b)}{|P_g(b)|}$$

a, b = anchors

$P_g(b)$ = all candidate pages of b

p_a, p_b = a candidate page of a, b

$\text{sim}(p_a, p_b)$ = relatedness between p_a & p_b

$P(p_b | b)$ = commonness of b in p_b

- Select top- ϵ pages based on the vote values for both anchors
- Select the best i in commoners, for both anchors

Result . . .

