

## Master on Foundations of Data Science



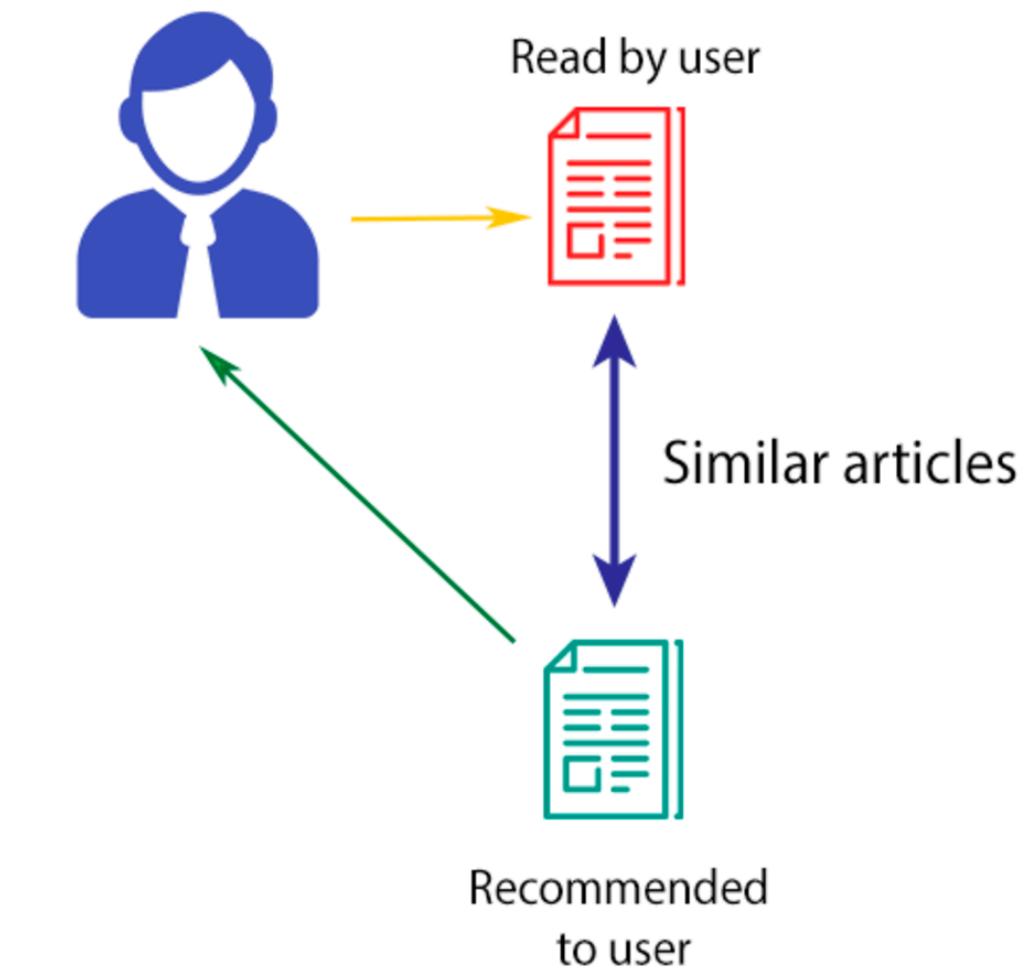
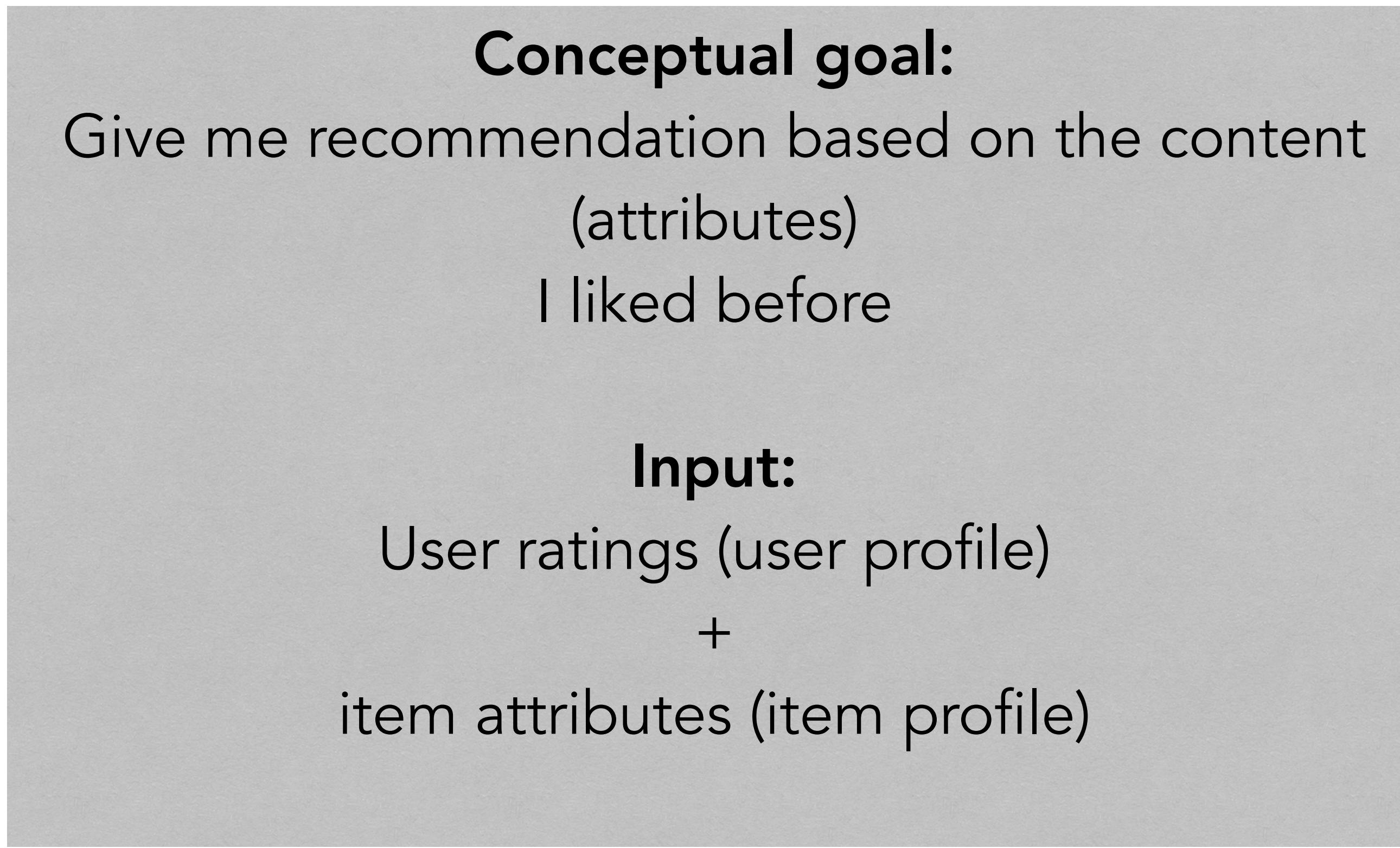
# Recommender Systems

Content Based Recomendations

Santi Seguí | 2022-2023

**Content-based** filtering assumes that a user will like items in the future that share features — like brand, cast, genre, etc. — with items they liked in the past

# Content-Based Methods



# Content-Based Filtering

- Requires **content** (from the items) that can be encoded as meaningful **features**.
  - Item title, description, price, image, etc...
- Need to compute a **similarity between items** based on the content of the items.
- **Users' tastes** must be represented as a **learnable function** of these content features.
- Does **not** exploit quality judgments of **other users**.
  - Unless these are somehow included in the content features.

# When Content Based?

Really popular for **cold-start** problems.

Popular in domain like:

**news** recommendation  
or **music** recommendation

# Advantages of CBRS

- **User independence**
  - CBRS exploit solely ratings provided by the active user to build the recommendation
  - No need for data on other users
- **Transparency**
  - Can provide explanations for recommended items by listing content-features that caused an item to be recommended
- **New Item (Cold Start on items)**
  - Can recommend new and unknown items

# Disadvantages

- **Limited to content.**
- **Over-specialization:** Content-based filtering provides a limited degree of novelty since it has to match up the features of a user's profile with available items.



## First rate two jokes.

Q: If a person who speaks three languages is called "trilingual," and a person who speaks two languages is called "bilingual," what do you call a person who only speaks one language?

A: American!

Less Funny

More Funny



Next

# Some Famous CB Recommender Systems

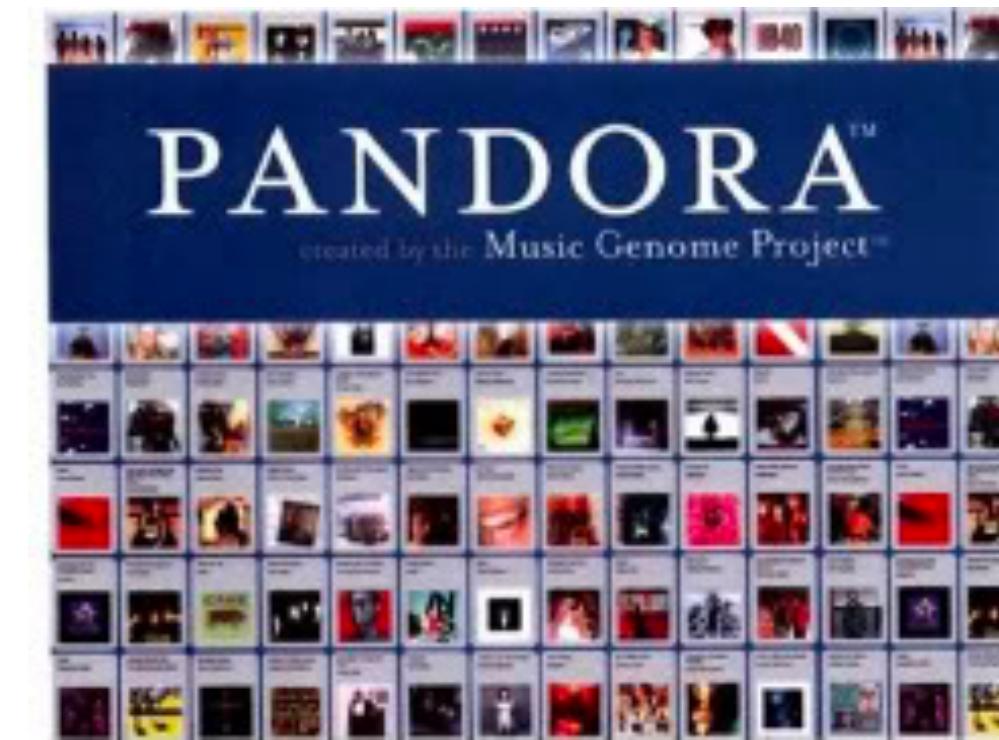


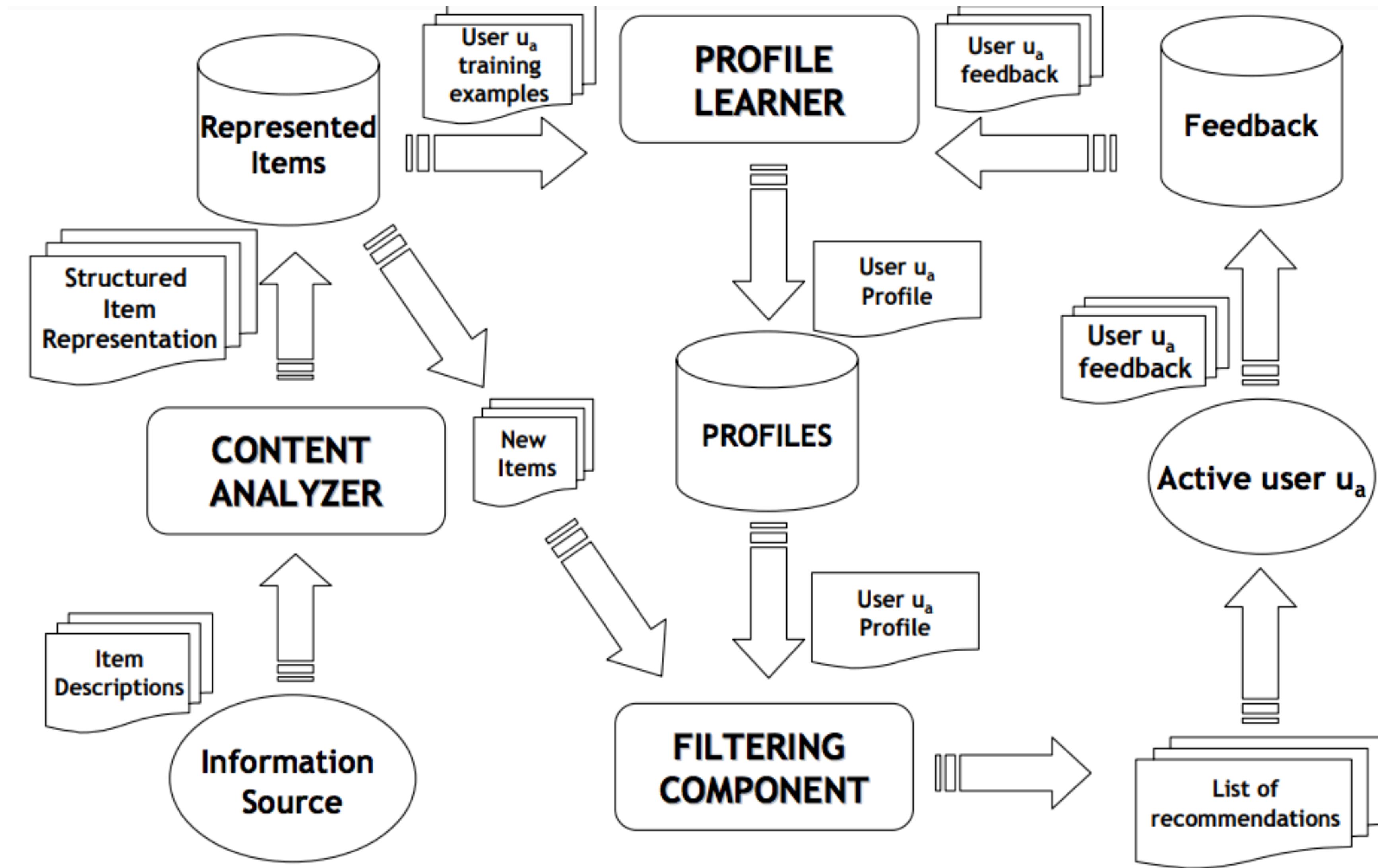
<https://www.pandora.com/>

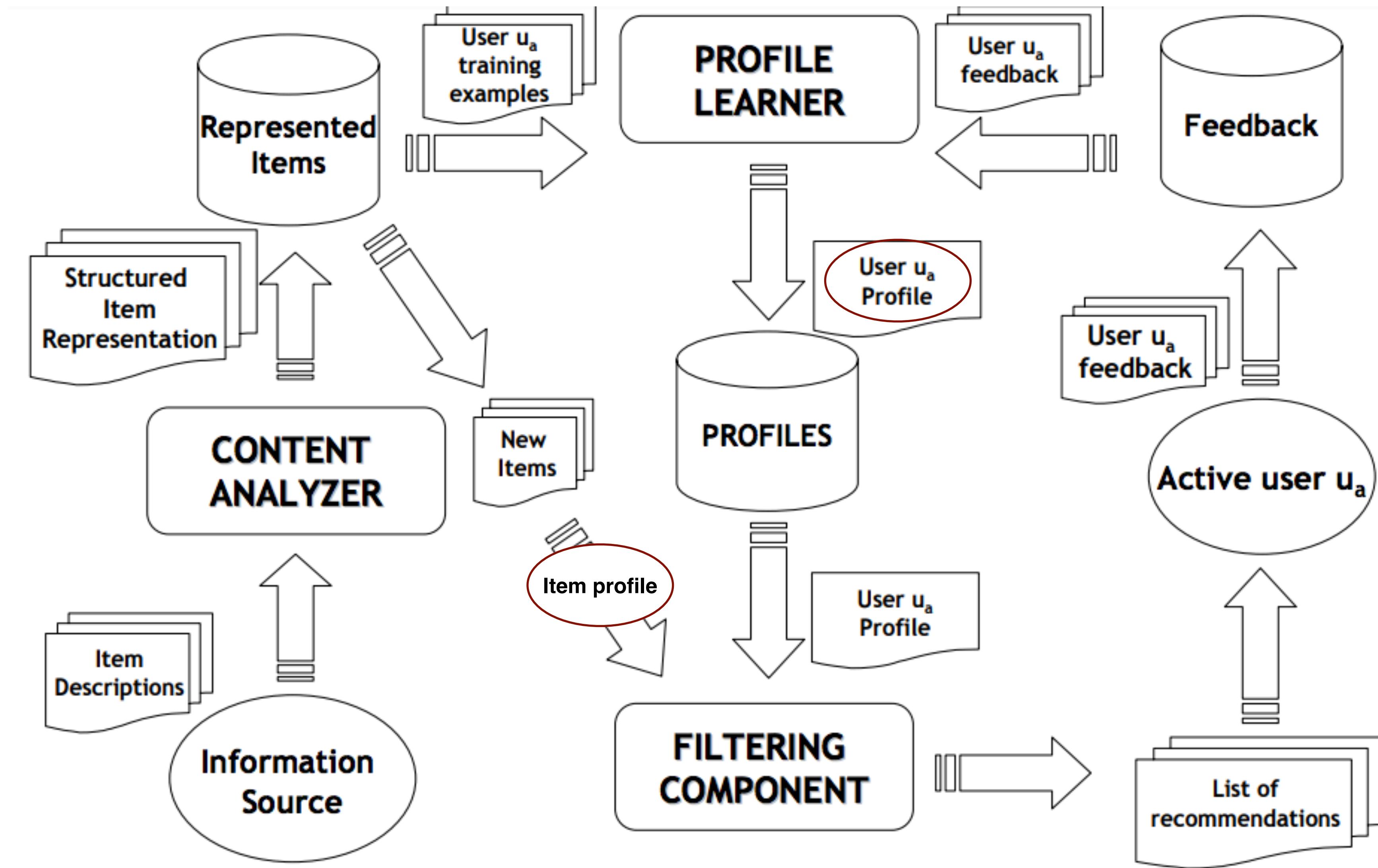


# Pandora

- How it works:
  - Base its recommendation on data from Music Genome Project
  - Assigns 400 attributes for each song, done by musicians.
  - Some reports say it takes half an hour per second of audio
  - Use this method to find songs which are similar to the user's favorite songs



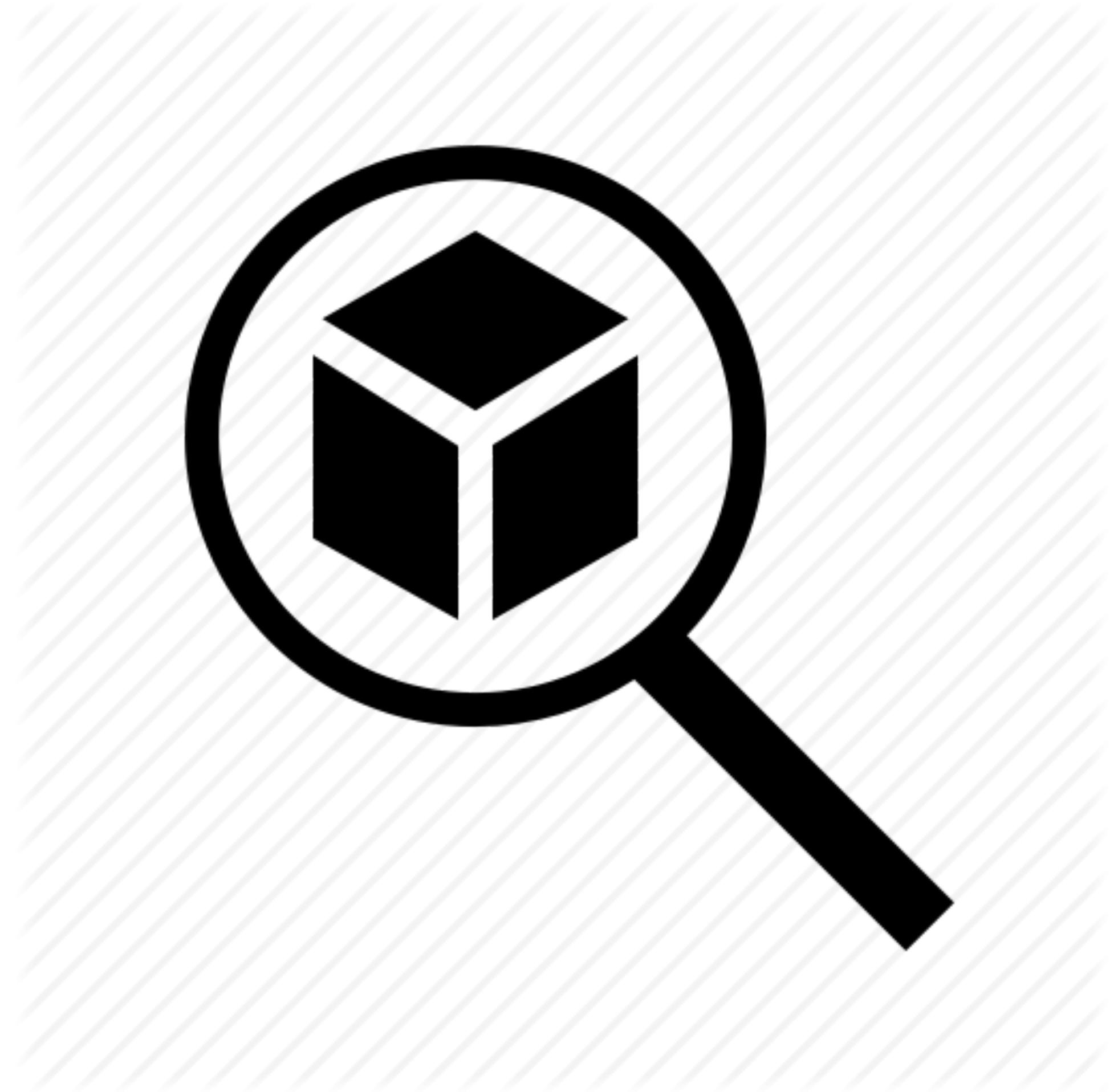




**“Key point”:**

**Similar Items** must have  
**similar** profile/vector **representation**

to do so, we need rich features and smart encoding



# Item Profile

# What is “content”

- Content Based recommenders systems have been applied mostly on **text document**
- Content from items, such as movies or songs, can be represented as text documents
  - With textual description of their basics characteristics
  - Structured: Each item is described with the same set of attributes
  - Unstructured: Free-text document

As for instance movies:



**Neruda (2016)** - [Limited]

R 107 min - Biography | Drama

Metascore: 88/100 (13 reviews)

An inspector hunts down Nobel Prize-winning Chilean poet, Pablo Neruda, who becomes a fugitive in his home country in the late 1940s for joining the Communist Party.

**Director:** Pablo Larraín  
**Stars:** Gael García Bernal, Luis Gnecco, Alfredo Castro, Pablo Derqui

[Watch Trailer](#) [Add to Watchlist](#)

# Item profile

- For each item, create an item profile
- Profile is a set of features.
  - Which features??
  - **Movies:** author, title, director, actor,...
  - **Images, videos:** raw content, metadata and tags
  - **People:** user profile, set of friends
  - **News:** keywords,..
- Convenient to consider the item profile as a vector:
  - One entry per feature (e.g., each actor, director, ..)
  - Vector might be boolean or real-valued



# We want to create a **Hotel Recommendation** system based on content

We have to create the item profile

Which features should we use?

# We want to create a **Hotel Recommendation** system based on content

from each hotel we have several features:

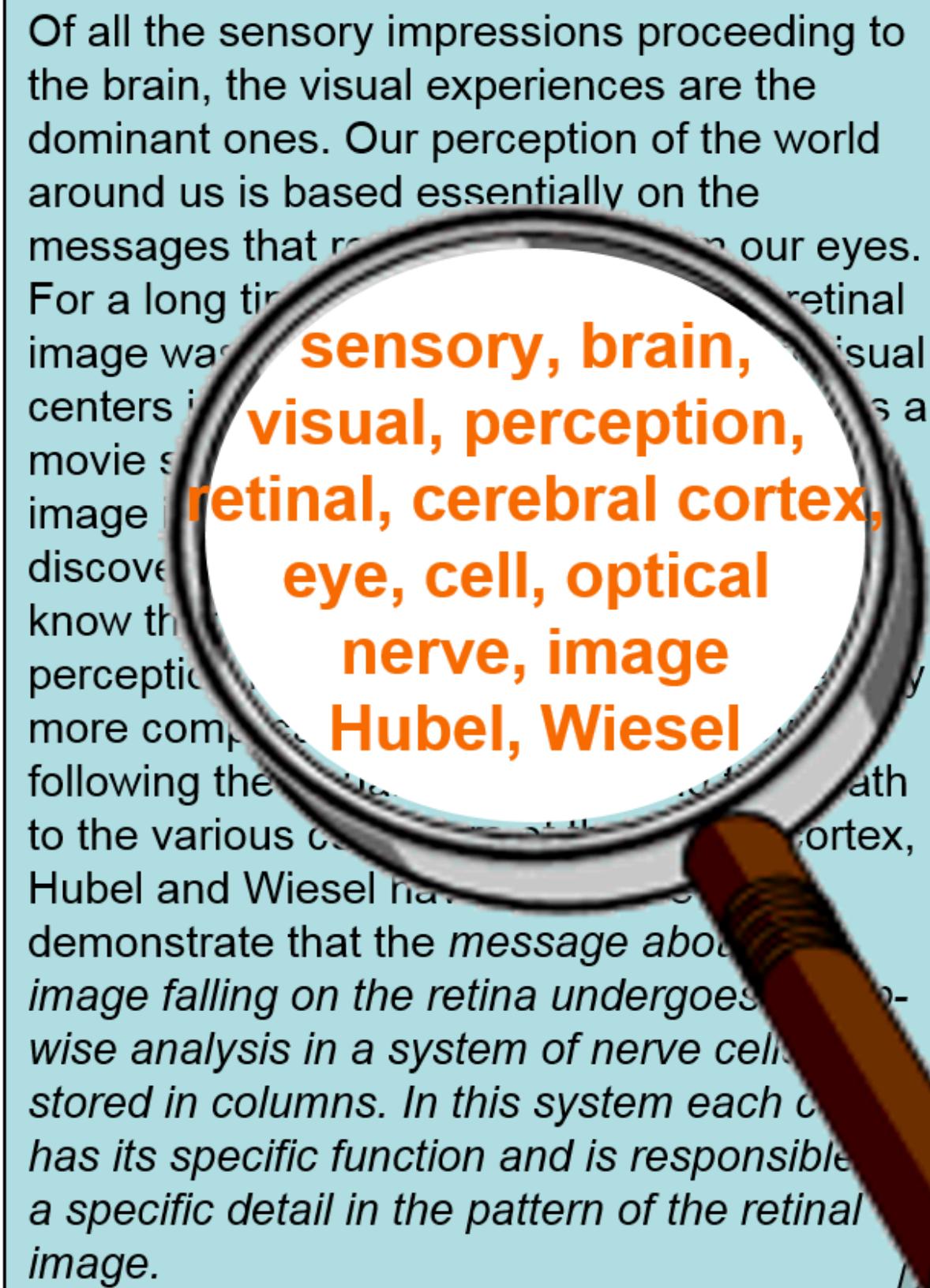
- City
- Location
- Price
- #Stars
- Swimming Pool?

How should be your feature vector?

And your user vector?

# How to describe textual information?

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our brain from our eyes. For a long time it was believed that the retinal image was processed directly in the visual centers in the cerebral cortex. It was a movie screen on which the image was projected. In 1960, two American scientists, David Hubel and Torsten Wiesel, discovered that the visual perception is much more complex than that. They found that the visual pathway follows the same basic pattern in all vertebrates. The image falling on the retina undergoes a top-down analysis in a system of nerve cells stored in columns. In this system each column has its specific function and is responsible for a specific detail in the pattern of the retinal image.



China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$560bn last year. The ministry also forecast imports would rise 20% to \$660bn. That would put China's trade balance at \$60bn, up from \$28bn in 2004. The US has been annoyed by the growth in China's trade surplus. The US says China's central bank, the People's Bank of China, deliberately keeps the value of the yuan low to help exports. The US says the Chinese government also needs to allow the yuan to rise more rapidly. The demand so far has been resisted by the Chinese government. The Chinese government has been allowed to trade within a narrow range of the dollar. But the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.



# Feature Representation and Cleaning

- Extremely important when the unstructured format is used for representation.
- Bag of Words (BOW) from the unstructured description of the products or Web Pages used to be used, however, these representations needs to be cleaned and represented in a suitable format for processing.
- Several important steps:
  - **Stop-word removal:** Words such “a”, “an,”, “the”, does not provide important information
  - **Stemming.** Variations of the same words are consolidated. For example, words such “hope” and “hoping” are consolidated into the common root “hop”
  - **Phrase extraction:** The idea is to detect words that occur together in documents on a frequent basis.

# TF-IDF

In information retrieval, **tf-idf**, short for term frequency-inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. It is often used as a weighting factor in information retrieval and text mining.

- **Term Frequency × Inverse Document Frequency**
- **Term Frequency** =
  - Number of occurrences of a term in a document (can be a simple count)
- **Inverse Document Frequency** =
  - Number of documents that contains this term
  - Typically :  $\text{Log}(\#\text{documents} / \#\text{documents with term})$

**So, items that appears rarely or appears everywhere are not important**

# TF-IDF

$$\text{tf}("this", d_1) = \frac{1}{5} = 0.2 \quad \text{tf}("this", d_2) = \frac{1}{7} \approx 0.14$$

Document 1		Document 2	
Term	Term Count	Term	Term Count
this	1	this	1
is	1	is	1
a	2	another	2
sample	1	example	3

idf is constant per corpus, and **accounts** for the ratio of documents that include the word "this". In this case, we have a corpus of two documents and all of them include the word "this".

$$\text{idf}("this", D) = \log\left(\frac{2}{2}\right) = 0$$

tf-idf is zero for the word "this", which implies that the word is not very informative as it appears in all documents.

$$\text{tfidf}("this", d_1) = 0.2 \times 0 = 0$$

$$\text{tfidf}("this", d_2) = 0.14 \times 0 = 0$$

# What does TFID do?

- Automatic find of stop words, common terms
- Promote core terms over accidental terms
- When it fails?
  - If core term is not used frequently in a document (e.g., legal contracts)

# Variants and Alternatives

- Some applications use variants on TF
  - Binary
  - Logarithmic frequencies
  - Normalized frequencies ( $\log(tf + 1)$ )

# Relevance and Problems

- Significance in Documents
  - Titles, heading,... (different weight?)
- Phrases and n-grams
  - “recommender system” != “recommender” and “system”
  - Adjacency
- General score
- Implied Content
  - Links, usage,...

# Keyword representation

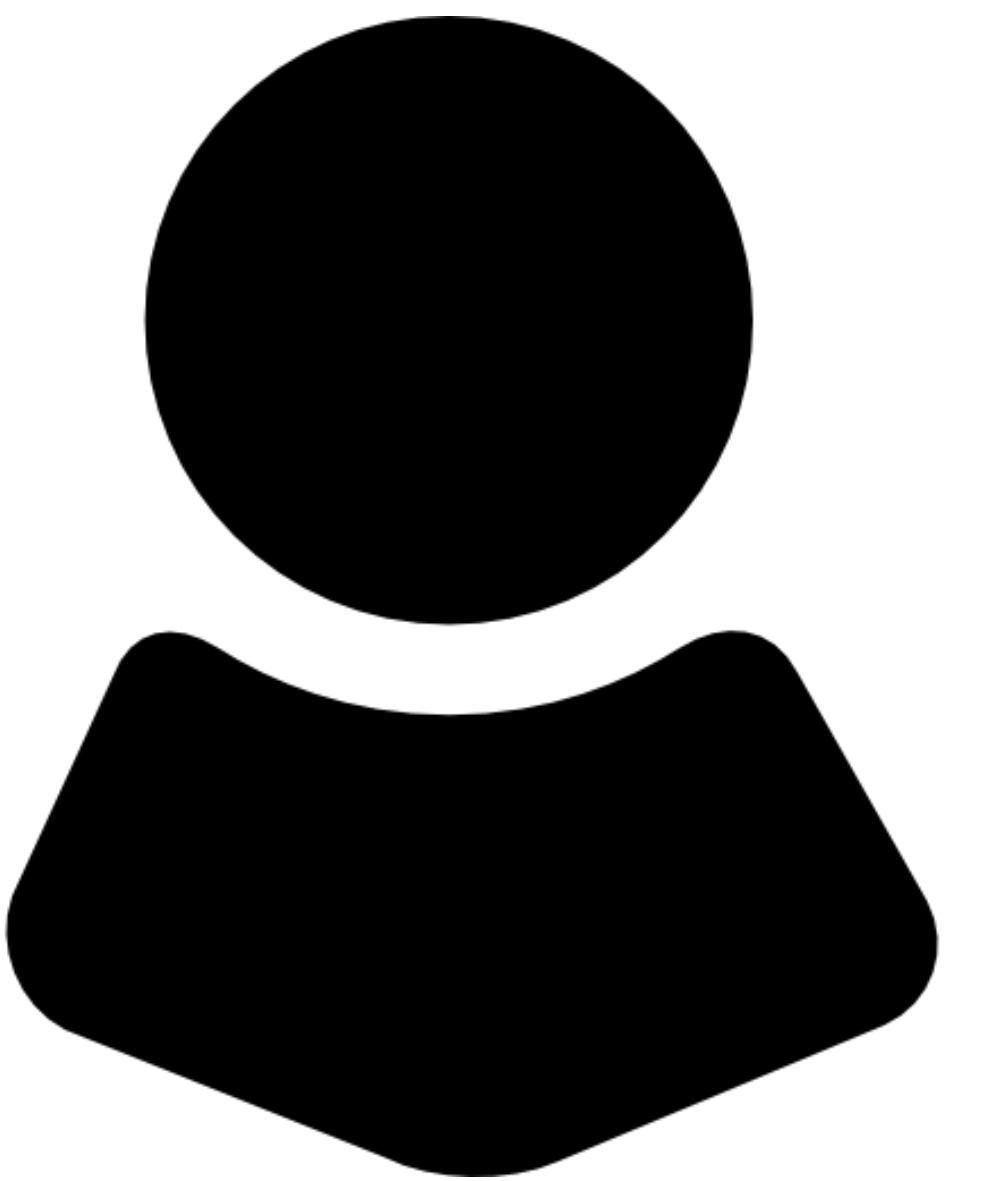
- The universe of possible keywords defines a content space
  - Each “keyword” is a dimension
  - Each item has a position in that space; that position defines a vector
  - Each user has a taste profile that is also a vector in that space
  - The match between user preferences and items is measured by how closely the two vectors align
  - May want to limit/collapse keyword space

# Vector Representation

- Simple 0/1 (keyword applies or does not)
- Simple occurrence count
- TF-IDF
- Other variants include factors such as document length
- Eventually, this vector is often normalized

# Other terms?

- Clothing attributes (color, size, etc..)
- Terms used in hotel reviews ( pool, front desk, child friendly)
- Terms used in news articles ( elections, football, economy)



User Profile

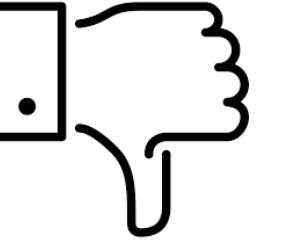
# User Preference

- My preferences:
  - **Movies** - I like SeVeN, American History X, Gladiator
  - **Hotels** - I prefer 24-hour front desk, internet, spa
  - **Music** - I like Blur, Pulp, The Verve,...

# User Preference

- **User Vector Space Model:**
  - single value for each dimension
  - same dimensions than **item Vector Space**

# How to build preferences?

- count the number of times the user chooses item with each keyword
- Set of “keywords” a user  or 
- or more sophisticated methods

# User Preferences

- How to **accumulate** features from the profiles?
  - Add together the item vectors?
  - Should we normalize first?
    - Should all items have the same weight?
    - Should we weight the vectors somehow?
      - We can use ratings..
      - Confidence?

# User Preferences

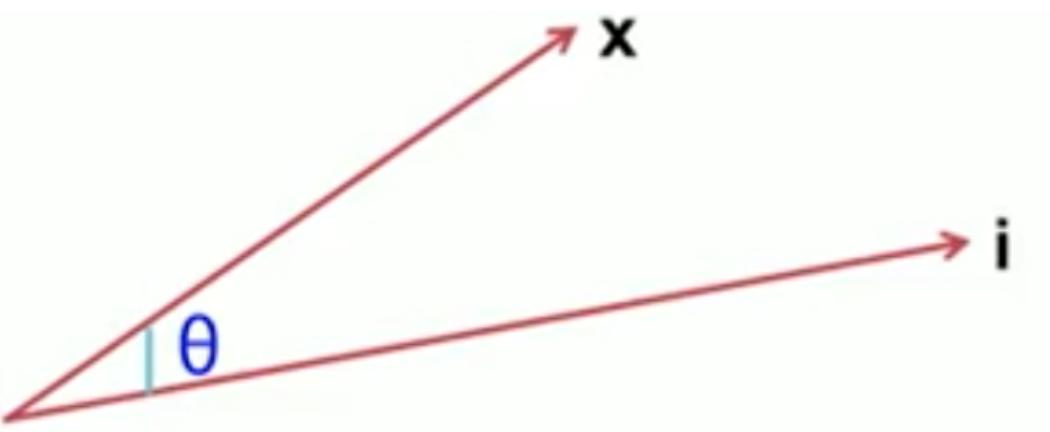
- What about **new user feedback from items?**
  - Update the user vector taking into account number of items used
  - Depending on the problem update with some decay



Making predictions

# Making Predictions

- User profile  $x$ , Item profile  $i$
- Estimate  $U(x, i) = \cos(\theta) = (x \cdot i) / (|x| |i|)$



- Technically, the cosine distance is actually the angle  $\theta$ . And the cosine similarity is the angle  $180 - \theta$

# How to improve it?

- Better classifier/Regressor
  - **Linear regression** to **XGboost** models are used
  - Each feature will have a different weight on the recommendation
- **Richer representations**



# ACM RecSys Challenge 2020



## Dataset description

The Data is available to download [here](#). Fields in each data entry are separated by the 1 character (*0x31 in UTF-8*) and each data entry will be characterized by the following features:

	Feature Name	Feature Type	Feature Description
Tweet Features	Text tokens Hashtags Tweet id Present media Present links Present domains Tweet type Language Timestamp	List[long] List[string] String List[String] List[string] List[string] String String Long	Ordered list of Bert ids corresponding to Bert tokeniza Tab separated list of hastags (identifiers) present in th Tweet identifier Tab separated list of media types. Media type can be i Tab separated list of links (identifiers) included in the t Tab separated list of domains included in the Tweet (t Tweet type, can be either Retweet, Quote, Reply, or Tc Identifier corresponding to the inferred language of th Unix timestamp, in sec of the creation time of the Twe
Engaged With User Features	User id Follower count Following count Is verified? Account creation time	String Long Long Bool Long	User identifier Number of followers of the user Number of accounts the user is following Is the account verified? Unix timestamp, in seconds, of the creation time of the
Engaging User Features	User id Follower count Following count Is verified? Account creation time	String Long Long Bool Long	User identifier Number of followers of the user Number of accounts the user is following Is the account verified? Unix timestamp, in seconds, of the creation time of the
Engagement Features	Engagee follows engager? Reply engagement timestamp Retweet engagement timestamp Retweet with comment engagement timestamp Like engagement timestamp	Bool Long Long Long Long	Does the account of the engaged tweet author follow If there is at least one, unix timestamp, in s, of one of t If there is one, unix timestamp, in s, of the retweet of t If there is at least one, unix timestamp, in s, of one of t If there is one, Unix timestamp, in s, of the like

How can we create a tweet similarity?

## Dataset description

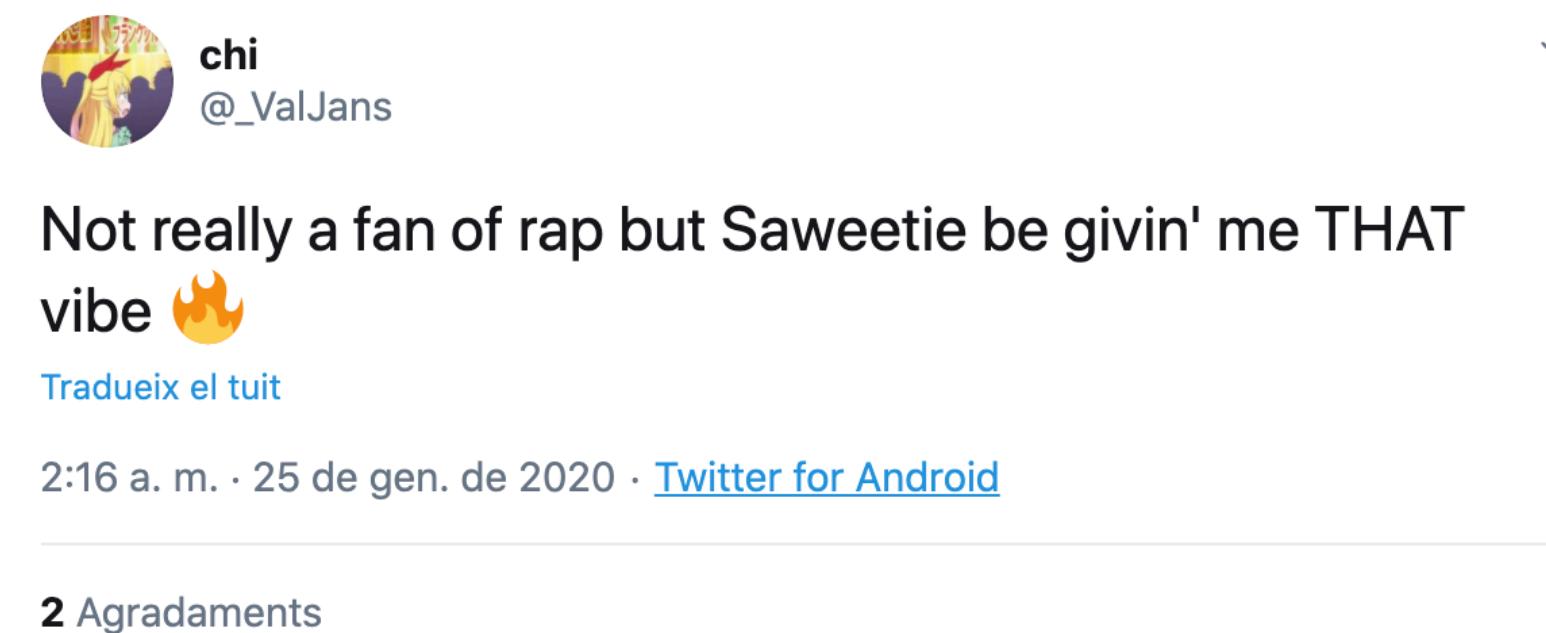
The Data is available to download [here](#). Fields in each data entry are separated by the 1 character (*0x31 in UTF-8*) and each data entry will be characterized by the following features:

	Feature Name	Feature Type	Feature Description
Tweet Features	Text tokens Hashtags Tweet id Present media Present links Present domains Tweet type Language Timestamp	<i>List[long]</i> <i>List[string]</i> <i>String</i> <i>List[String]</i> <i>List[string]</i> <i>List[string]</i> <i>String</i> <i>String</i> <i>Long</i>	Ordered list of Bert ids corresponding to Bert tokeniza Tab separated list of hastags (identifiers) present in th Tweet identifier Tab separated list of media types. Media type can be i Tab separated list of links (identifiers) included in the Tab separated list of domains included in the Tweet (t Tweet type, can be either Retweet, Quote, Reply, or Tc Identifier corresponding to the inferred language of th Unix timestamp, in sec of the creation time of the Twe
Engaged With User Features	User id Follower count Following count Is verified? Account creation time	<i>String</i> <i>Long</i> <i>Long</i> <i>Bool</i> <i>Long</i>	User identifier Number of followers of the user Number of accounts the user is following Is the account verified? Unix timestamp, in seconds, of the creation time of the
Engaging User Features	User id Follower count Following count Is verified? Account creation time	<i>String</i> <i>Long</i> <i>Long</i> <i>Bool</i> <i>Long</i>	User identifier Number of followers of the user Number of accounts the user is following Is the account verified? Unix timestamp, in seconds, of the creation time of the
Engagement Features	Engagee follows engager? Reply engagement timestamp Retweet engagement timestamp Retweet with comment engagement timestamp Like engagement timestamp	<i>Bool</i> <i>Long</i> <i>Long</i> <i>Long</i> <i>Long</i>	Does the account of the engaged tweet author follow If there is at least one, unix timestamp, in s, of one of t If there is one, unix timestamp, in s, of the retweet of t If there is at least one, unix timestamp, in s, of one of t If there is one, Unix timestamp, in s, of the like

How can we create a tweet similarity?

# Bert

**BERT** (Bidirectional Encoder Representations from Transformers) is a recent paper published by researchers at Google AI Language. It has caused a stir in the Machine Learning community by presenting state-of-the-art results in a wide variety of NLP tasks, including Question Answering (SQuAD v1.1), Natural Language Inference (MNLI), and others.



chi  
 @\_ValJans

Not really a fan of rap but Saweetie be givin' me THAT  
vibe 🔥

Tradueix el tuit

2:16 a. m. · 25 de gen. de 2020 · [Twitter for Android](#)

2 Agradaments

## Text is codified into tokens

```
"[CLS] Not really a fan of rap but Saweetie be givin'me THAT vibe [UNK] [SEP]"
```

```
'101 16040 30181 169 10862 10108 35562 10473 74666 23203 10400 10347 38356 15478 112 10911 157 58132 11090 13956 11044 100 102'
```

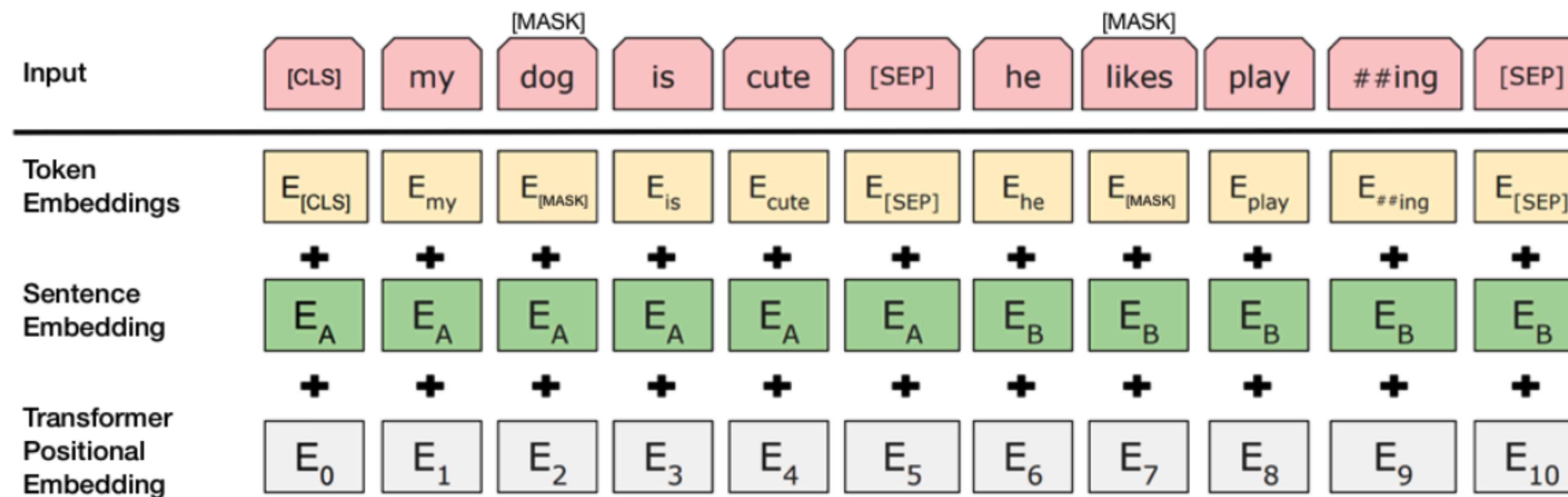
How can we create a tweet similarity?

# Bert

## BERT (language model)

From Wikipedia, the free encyclopedia

**Bidirectional Encoder Representations from Transformers (BERT)** is a technique for NLP (Natural Language Processing) pre-training developed by Google. BERT was created and published in 2018 by Jacob Devlin and his colleagues from Google.<sup>[1][2]</sup> Google is leveraging BERT to better understand user searches.<sup>[3]</sup>



<https://github.com/google-research/bert>

'101 10808 10173 61644 12387 11132 13474 10108 31206 37715 10251 117 10462 10571 32719 119 14120 131 120 120 188 119 11170 120  
12428 63051 11537 10138 11369 11373 11259 10477 14120 131 120 120 188 119 11170 120 170 11396 11779 54889 14703 10729 11281 11166  
10731 102'

US confirms second case of coronavirus, 50 under investigation. <https://t.co/76krNuL9Ev> <https://t.co/b8VGCY2I5S> [SEP]



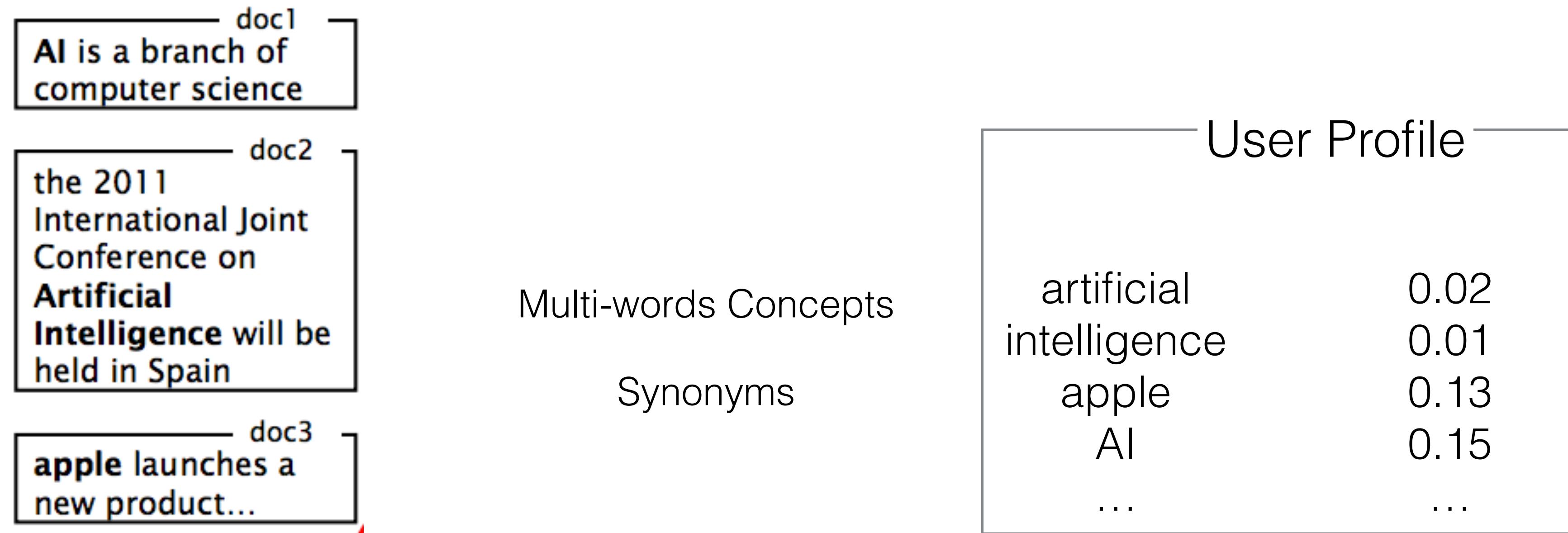
## Some similar tweets?

(1.0000002, '[CLS] US confirms second case of coronavirus, 50 under investigation. <https://t.co/76krNuL9Ev> <https://t.co/b8VGCY2I5S> [SEP]')  
(0.8663348, '[CLS] China coronavirus : Death toll rises as disease spreads <https://t.co/Cko9Z1fTRZ> [SEP]')  
(0.8271704, '[CLS] Las autoridades chinas confirman el primer caso de curación del nuevo coronavirus <https://t.co/yYldx6ATYR> [SEP]')  
(0.8174852, "[CLS] Inside President Trump's high - stakes impeachment defense effort <https://t.co/qWHQ8mAfQW> <https://t.co/FJ1a330df6> [SEP]")  
(0.7904908, '[CLS] NEW PEER NEWS PIECE!.. Person of the Year : The French Worker on Strike. <https://t.co/EGaxAvhh7D> <https://t.co/zGkMSnQPtm> [SEP]')  
(0.7811252, '[CLS] AB6IX make dumplings & ; tteokguk by hand for Lunar New Year!. <https://t.co/hYGrYmHBK1> <https://t.co/48t5Rbz4uF> [SEP]')  
(0.76791394, "[CLS] Alaska health officials monitoring Coronavirus, deferring to CDC's guidance on when additional measures may be needed : <https://t.co/L9SLo2Ezwi> [SEP]")  
-----

# How to improve content encoding?

Keywords are **not appropriate** for representing content, due to **polysemy, synonymy, multi-word concepts**,....

# Keyword-based Models



apple  
Polysemy



**NLP methods are needed for the elicitation of user interests**

# Richer representations

- Semantic Analysis
  - **Semantics:** concept identification in text-based representations through advanced NLP techniques -> “beyond keywords”
  - **Personalization:** representation of user information needs in an effective way -> “deep user profiles”

# Matrix Factorization

- Latent Semantic Analysis
- Latent Dirichlet Allocation

## LSA

$$T = K \times S \times D^T$$

Document by Keyword Matrix ( $d \times k$ )

Topic by Keyword Matrix ( $z \times k$ )

Topic by Topic Matrix ( $z \times z$ )

Document by Topic Matrix ( $d \times z$ )

## LDA

$$P(k|d) = P(k|z) \times P(z|d)$$

Document distribution over Keywords ( $d \times k$ )

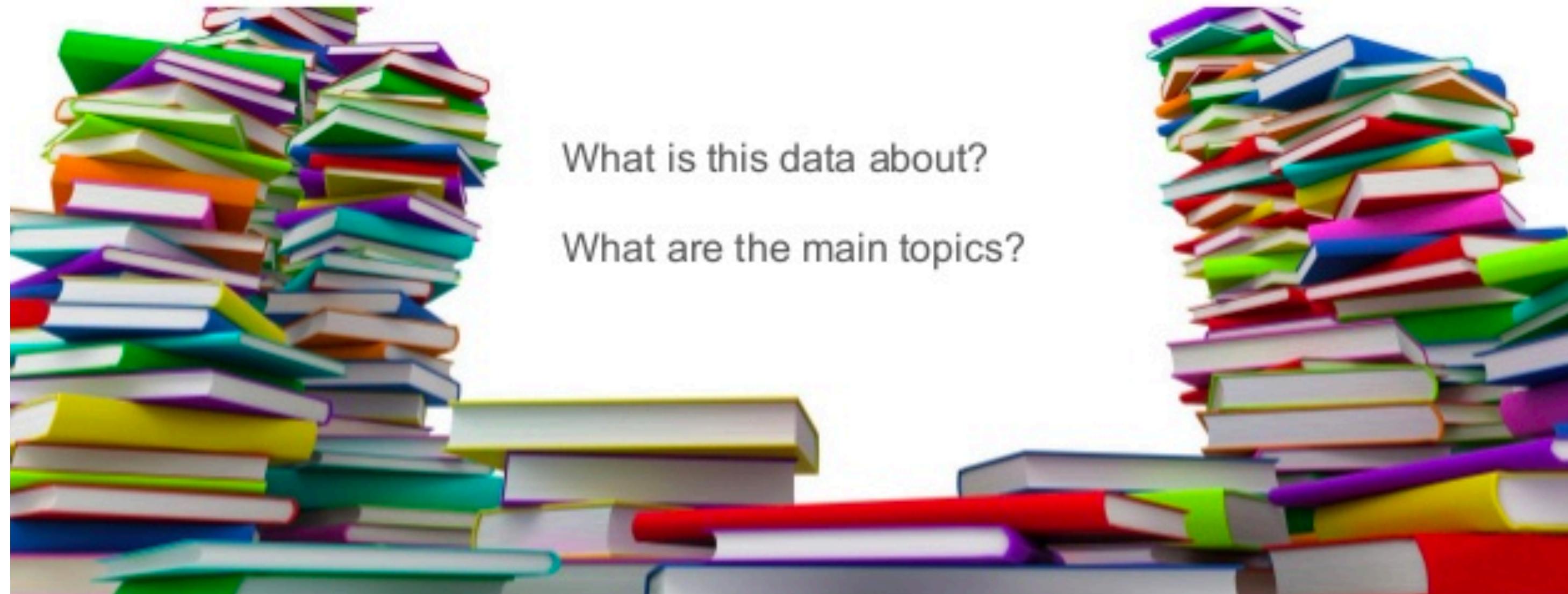
Topic distribution over Keywords ( $z \times k$ )

Document distribution over Topics ( $d \times z$ )

Fig. 2: Matrix decomposition for LSA and LDA.

# Topic Modeling

A simple way to analyze topics of large text collections (corpus).



What is this data about?

What are the main topics?

Latent dirichlet allocation

DM Blei, AY Ng, MI Jordan

Journal of machine Learning research 3 (Jan), 993-1022

17706 2003

Topic 1		Topic 2		Topic 3	
term	weight	term	weight	term	weight
game	0.014	space	0.021	drive	0.021
team	0.011	nasa	0.006	card	0.015
hockey	0.009	earth	0.006	system	0.013
play	0.008	henry	0.005	scsi	0.012
games	0.007	launch	0.004	hard	0.011

Latent dirichlet allocation

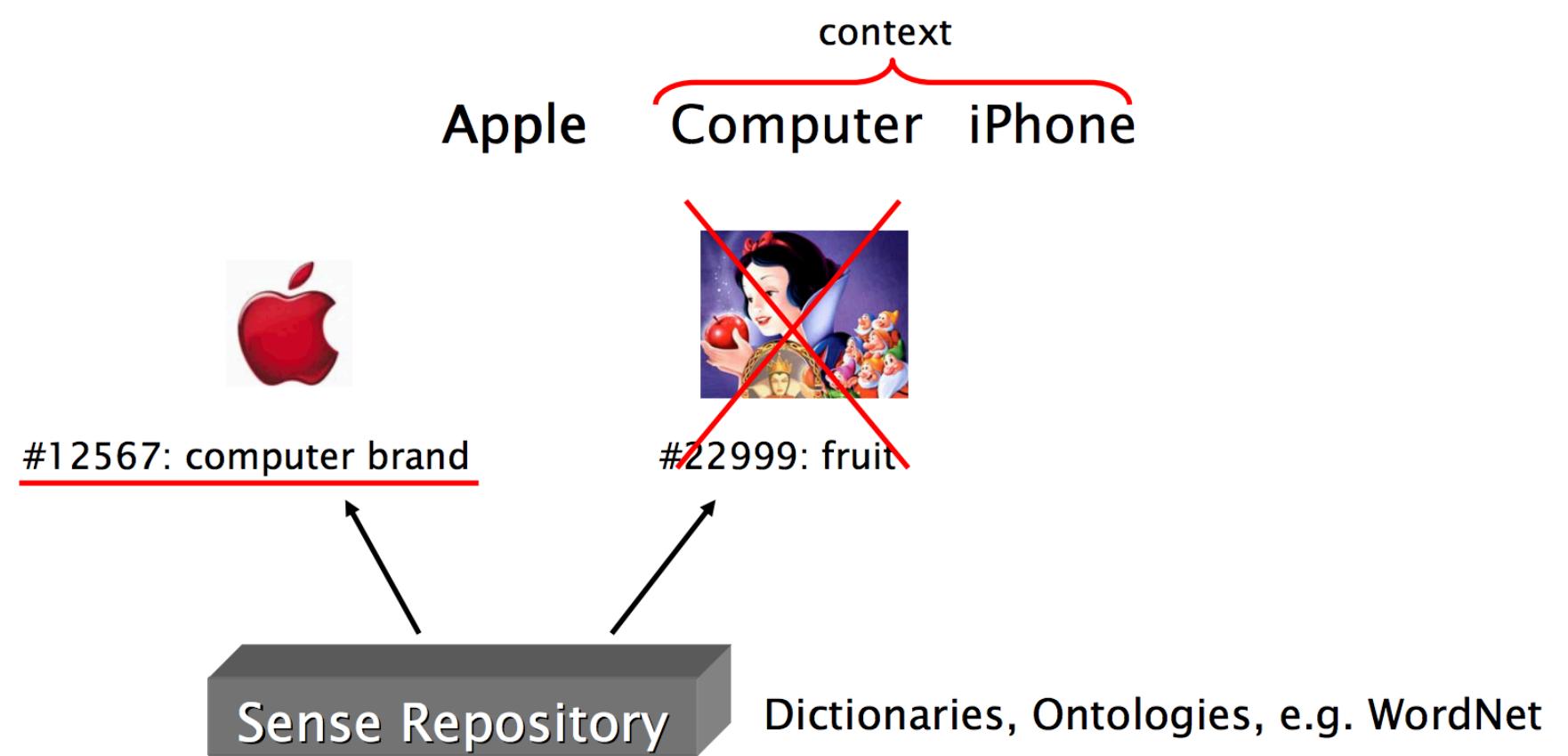
DM Blei, AY Ng, MI Jordan

Journal of machine Learning research 3 (Jan), 993-1022

17706 2003

# Semantic Analysis using Ontologies

- Word sense Disambiguation (WSD) -> From words to meanings
  - WSD selects the proper meaning (sense) for a word in a text by taking into account the context in which that word occurs



A content-collaborative recommender that exploits WordNet-based user profiles for neighborhood formation

M Degenni, P Lops, G Semeraro

User Modeling and User-Adapted Interaction 17 (3), 217-255

181

2007

# Semantic Analysis using Encyclopedic Knowledge Sources

- Word2Vect (you will see it in DeepLearning)

Wikipedia is viewed as an **ontology** - a collection of ~**1M** concepts

Every Wikipedia article represents a **concept**

Panthera

From Wikipedia, the free encyclopedia

**Panthera** is a genus of the family Felidae (the cats) which contains four well-known living species: the lion, tiger, jaguar, and leopard. The genus comprises about half of the big cats. One meaning of the word *panther* is to designate cats of this family. Only these four cat species have the anatomical changes enabling them to roar. The primary reason for this was assumed to be the incomplete ossification of the hyoid bone. However, new studies show that the ability to roar is due to other morphological features, especially of the larynx. The snow leopard (*Uncia uncia*, which is sometimes included within *Panthera*, does not roar. Although it has an incomplete ossification of the hyoid bone, it lacks the special morphology of the larynx, which is typical for lions, tigers, jaguars and leopards.<sup>[1]</sup>

Species and subspecies



Tiger

Scientific classification

Kingdom: Animalia

Phylum: Chordata



Panthera

Cat [0.92]

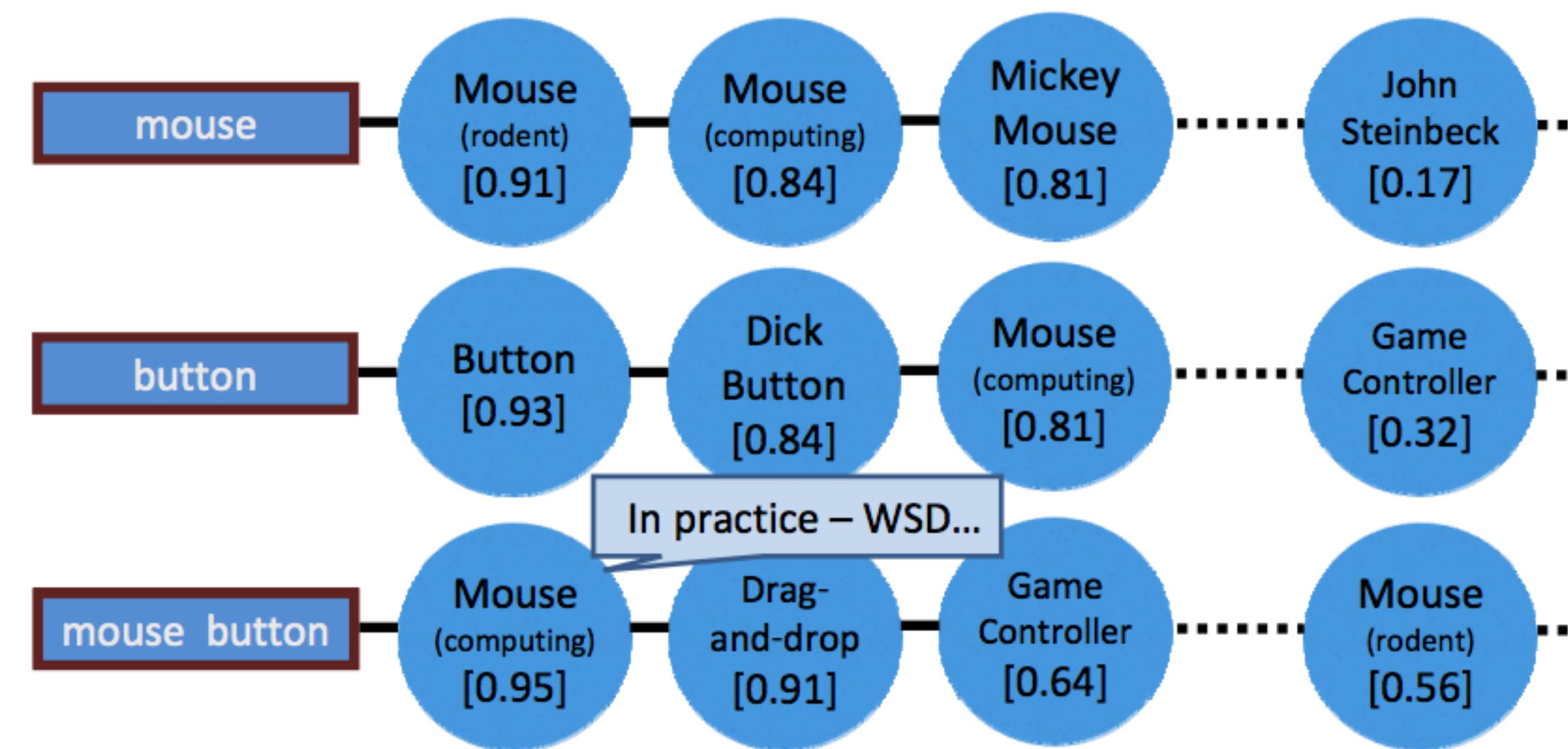
Leopard [0.84]

Roar [0.77]

Article words are associated with the concept (TF-IDF)

# Semantic Analysis using Encyclopedic Knowledge Sources

The **semantics** of a **text fragment** is the **average vector (centroid)** of the **semantics of its words**



# word2vec

(WATER - WET) + FIRE = FLAMES

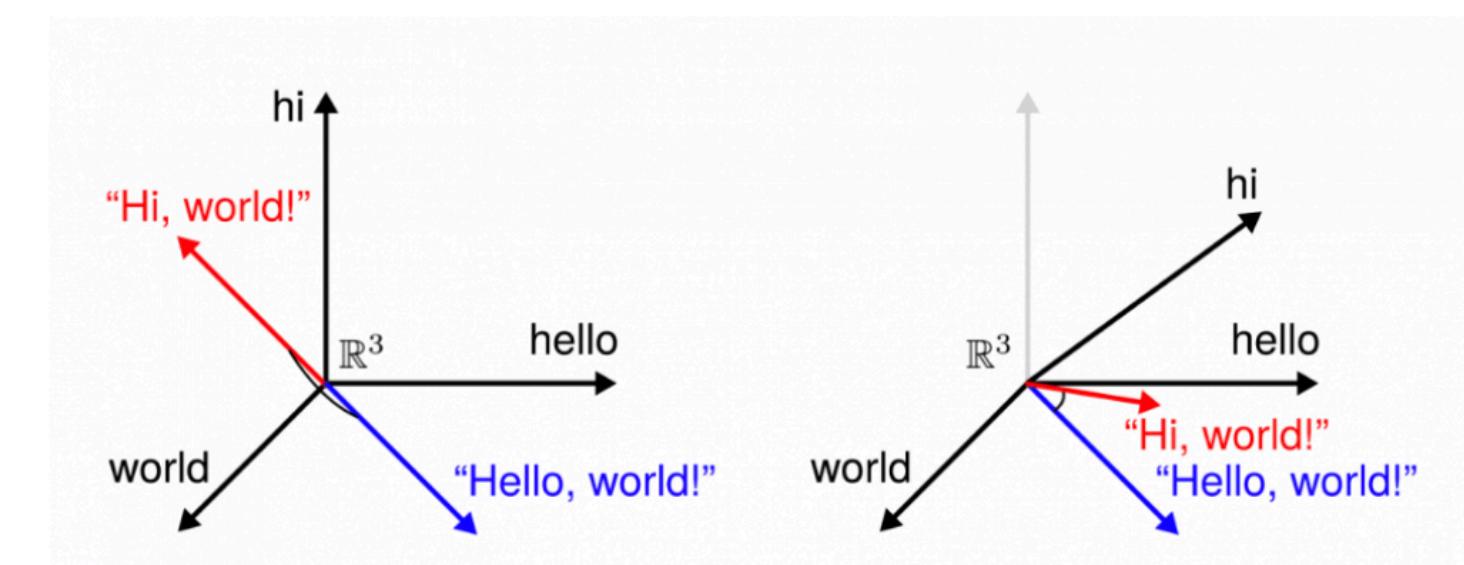
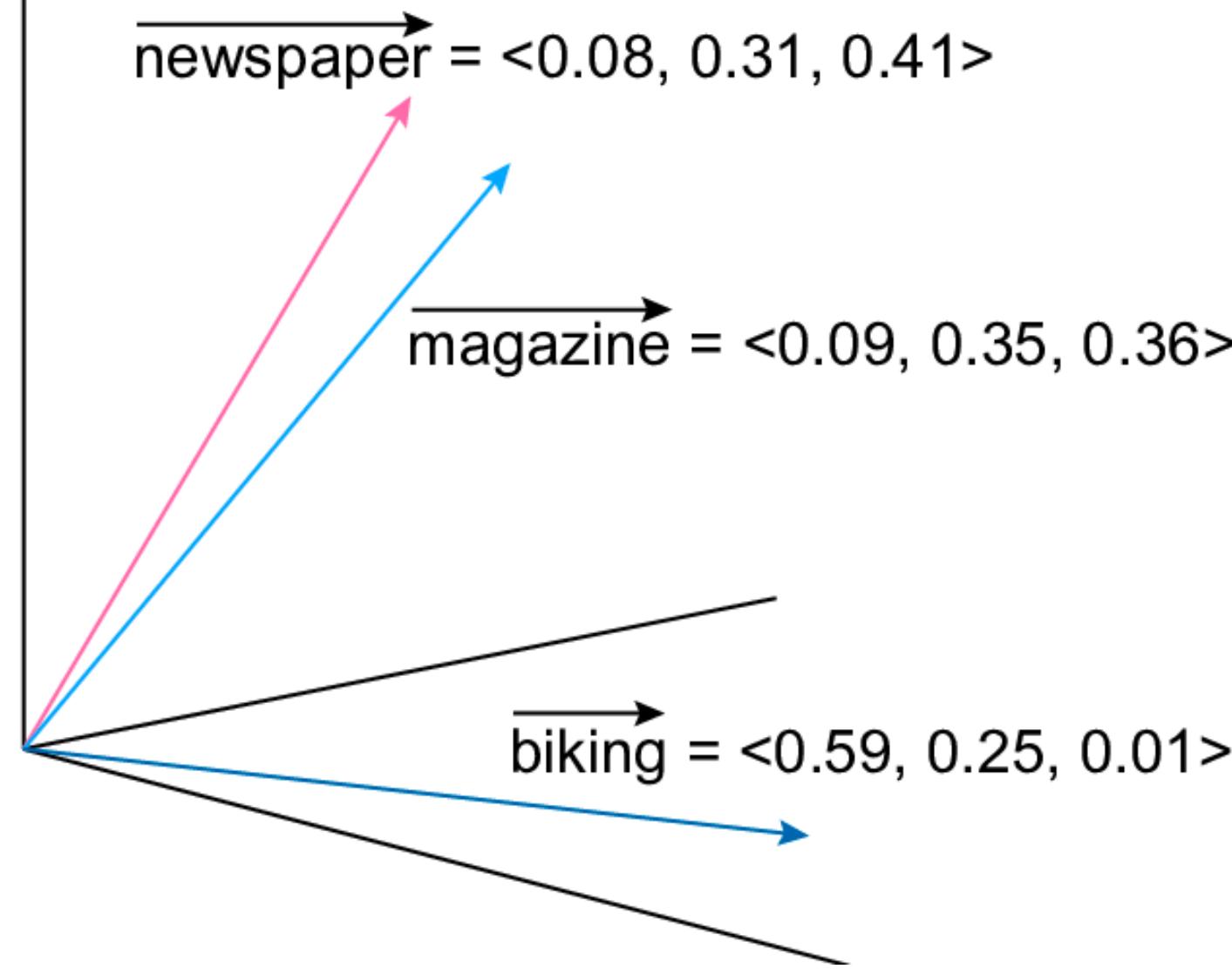
(PARIS - FRANCE) + ITALY = ROME

(WINTER - COLD) + SUMMER = WARM

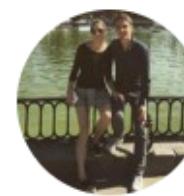
(MINOTAUR - MAZE) + DRAGON = SIMCITY

- █ : Target Word
- █ : Context Word

- c=0    The cute cat jumps over the lazy dog.
- c=1    The cute cat jumps over the lazy dog.
- c=2    The cute cat jumps over the lazy dog.



# Text Similarities : Estimate the degree of similarity between two texts



Adrien Sieg [Follow](#)

Jul 5, 2018 · 29 min read ★



*Note to the reader: Python code is shared at the end*

We always need to **compute the similarity in meaning between texts**.

- **Search engines** need to model the relevance of a document to a query, beyond the overlap in words between the two. For instance, **question-and-answer** sites such as Quora or Stackoverflow need to determine whether a question has already been asked before.
- In **legal matters**, text similarity task allow to mitigate risks on a new

... 1 2 3 4 ...

<https://medium.com/@adriensieg/text-similarities-da019229c894>