# The Human Factor in Big Data Analysis

Karl Aberer, EPFL
Distributed Information Systems Laboratory
lsir.epfl.ch

Distributed Information Systems Laboratory

Big Data = Volume, Velocity, Variety, Veracity

**Variety = Semantics = Meaning**
Retrieval, data integration, information extraction, ...

**Veracity = Pragmatics = Utility**
Data quality, credibility, authority, trust, ...

Stating the obvious: every semantic and pragmatic information processing task related to human concerns requires human input

For example: Google is a huge relevance feedback engine

Distributed Information Systems Laboratory

Big Data analysis today

- **Key innovation:** capacity to automatically process and analyse huge volumes of data

- **Key bottleneck:** human input to make the processing meaningful

Example: recent progress in machine translation and image recognition with deep learning

- Rely on huge corpuses with "ground truth"
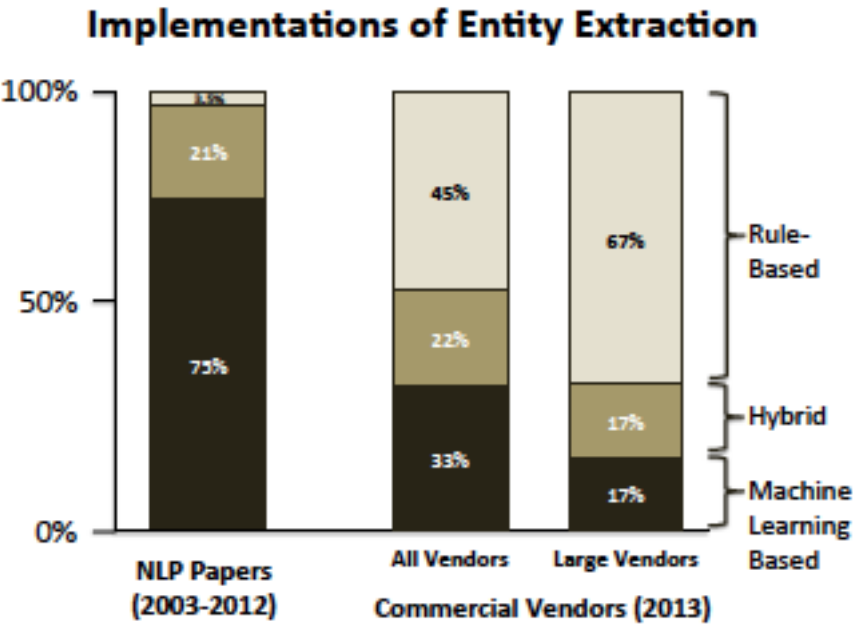
Distributed Information Systems Laboratory



Vinyals, Oriol, et al. "Show and tell: A neural image caption generator."
*Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015.

4

Distributed Information Systems Laboratory

CHRIS ANDERSON    MAGAZINE    06.23.08    12:00 PM

# THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE

**No models!**
**No causality!**
**No understanding!**

- But often no ground truth available, in particular for applications with "not so big data" and involving expert knowledge

**Implementations of Entity Extraction**



Chiticariu, Laura, Yunyao Li, and Frederick R. Reiss. "Rule-Based Information Extraction is Dead! Long Live Rule-Based Information Extraction Systems!." *EMNLP*. No. October. 2013.

Distributed Information Systems Laboratory

## Three case studies

1. Web credibility

   *How human input enables machine learning*
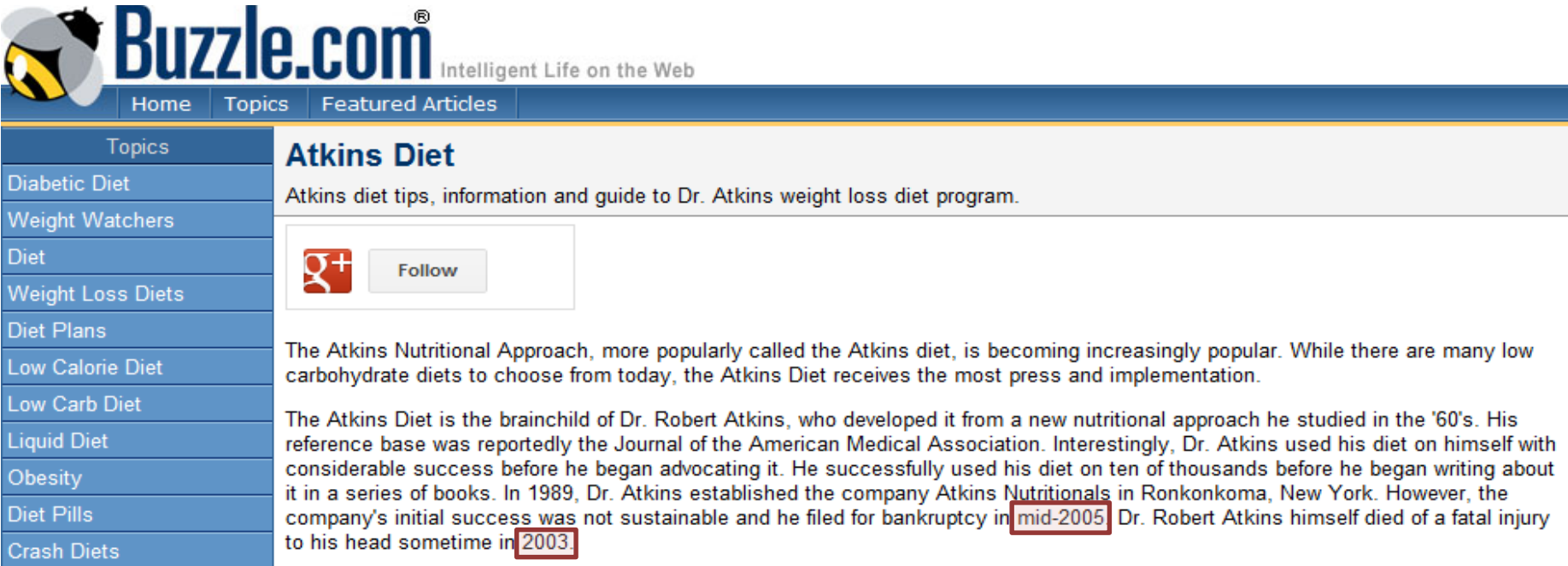
2. Data integration

   *How humans and machines cooperate efficiently in a problem solving task*

3. Social Media Analytics

   *How humans interpret latent structures found by machine learning*

# CASE STUDY 1: WEB CREDIBILITY

# How to Evaluate Web Credibility?



*When you browse a webpage, how do you know it's content is valid and accurate?*

# Web Credibility – the problem

Increasingly difficult to assess credibility of Web content

- Economic incentives to manipulate information
  - Marketing, fraud, political motives, etc.
- Enormous volume of web information

User: Believe or not?

"looks" benign

Presentation features

non-credible site

Adversary: Put $$$ to make it look credible?

# Which features indicate credibility?

- Numerous candidate features could indicate credibility?

- How to determine?

- Let experts annotate a collection of documents

**École Polytechnique Fédérale de Lausanne**

Distributed Information Systems Laboratory

| Topic | Query Terms | Expert URL Filters | # of Users |
|---|---|---|---|
| Health | Atkins diet effectiveness<br>P90x exercise program<br>H1N1 vaccine side effects<br>Alzheimer's genes<br>Autism warning signs | ncbi.nlm.nih.gov/pubmed<br>pubmedcentral.nih.gov | 254,175 |
| Finance | Is it a good time to invest in gold?<br>What mutual funds to invest in<br>Reduce personal debt<br>Mortgage refinancing<br>Is it a good time to invest? | bloomberg.com<br>edgar-online.com<br>hoovers.com<br>sec.gov | 201,014 |
| Politics | Iran election rigged<br>Cash for clunkers eligibility<br>Obama birthplace<br>Death Panels<br>Tea Party | foreignaffairs.com<br>theatlantic.com<br>foreignpolicy.com<br>hir.harvard.edu<br>economist.com | 66,155 |
| Celebrity News | Lady Gaga<br>Adam Lambert<br>Nadya Suleman<br>Floyd Landis<br>Michael Jackson | ew.com<br>usmagazine.com<br>people.com | 692,611 |
| Environmental Science | Renewable energy<br>Green jobs<br>Climate change<br>Cap-and-trade<br>Organic Eating | pewclimate.org<br>epa.gov<br>rff.org<br>nrdc.org<br>whitehouse.gov/administration/ceq | 83,476 |
| *All Users* | | *(none)* | 50,473,520 |

- Corpus of 1000 documents
- Evaluated by domain experts
- (prepared by MS Research)

Statistical tests

- Identification of features providing the signals on credibility

**Informativeness**
**Google Search Ranking**
Domain Type (.gov, .edu)
Popularity on Twitter
**Readability**
Number of Bookmarks
Ads Prominence    Objectivity
**Use of Punctuation**
Use of Grammar
**Webpage Design**
Web Graph Structure
Popularity on Facebook
Browsing Patterns
Text Complexity
**Webpage Topic**

The content of a webpage as well as the social popularity offer signals for credibility

Distributed Information Systems Laboratory

**EPFL**
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE
Distributed Information Systems Laboratory

**generation** rescue
*hope for recovery*

enter keywords

Home | Recovery | Newly Diagnosed | Resources | Blog | Events | Store | Donate

◄ **Recovery** » ◄ **Biomedical Treatment** » Treatments to Explore

**Biomedical Treatment**

Three Steps

The Big Picture

Treatments to Explore

# Treatments to Explore

Conventional medicine treats the symptoms of autism. Biomedical treatment addresses the root cause.

There is a wealth of biomedical therapies that treat the underlying issues of autism inside the body.

The following is a list of biomedical treatments to explore with a physician in order to help heal the body:

**1. Follow the gluten-free, casein-free, soy-free diet and remove other food allergens.**

The yeast-autism connection can be a result of **candida** (type of yeast) overgrowth in the system. This leads to many different behaviors such as, fogginess, sensory issues, negative behaviors.

"Gluten and Dairy seem to affect a lot of our children with autism and thus we see a

**Donate**
Support the Cause.

**Get Started**
Steps to Support and Recovery.

**Biomedical Grant**
Apply for a Grant Online.

## Events

**Feb. 20, 2014**
#ChatAutism with Dr. Bo Wagner

**Feb. 22 - 23, 2014**
Gluten Free Allergy Free Expo Phoenix

**Mar. 20, 2014**
#ChatAutism with Dr. Douglas Bibus

**Mar. 27 - 29, 2014**
MAPS Spring 2014 Clinician CME Training Conference

# Not so credible statements

**1. Follow the gluten-free, casein-free, soy-free diet and remove other food allergens.**

The yeast-autism connection can be a result of **candida** (type of yeast) overgrowth in the system. This leads to many different behaviors such as, fogginess, sensory issues, negative behaviors.

"Gluten and Dairy seem to affect a lot of our children with autism and thus we see a lot of children respond terrifically when these are removed from the diet. The goal behind changing diets is to remove chemicals, toxins and potential neurotransmitters, which are liberated when food are broken down. These substances could be toxic for the brain and cause behavioral trouble in kids who are sensitive. Whether kids test as allergic or not, often they are causing a negative effect on the child and they must be removed. Each child has his or her own set of sensitivities that he or she can't deal with properly. When we change their diets, 80 percent of the kids with autism seem to respond." - Dr. Jerry Kartzinel, from "**Healing and Preventing Autism**" by Jenny McCarthy and Dr. Jerry Kartzinel.

- Effectiveness of the **gluten-free, casein-free diet** for children diagnosed with autism spectrum disorder: based on parental report.

- **Nutrition Guide** and how to implement the GFCF diet;

- Dr. Jerry's blog on why to implement the gluten free, casein free diet - **Parts 1** & **Part 2**

**More Resources:**

- **GFCFDiet.com**
- **The role of Clostridia and Autism**
- **The Yeast Problem and Bacteria By-products**
- **Improved Diet Helps Children with Autism**

**Feb. 22 - 23, 2014**
Gluten Free Allergy Free Expo
Phoenix

**Mar. 20, 2014**
#ChatAutism with Dr. Douglas Bibus

**Mar. 27 - 29, 2014**
MAPS Spring 2014 Clinician CME
Training Conference

Sponsors

Supporting ... you

ENZYMEDICA
The Enzyme Experts

ANGELICA

blk.

intelliBED
Perfecting the Science of Sleep

Oxy HEALTH
Portable Hyperbarics

16

**EPFL**
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Distributed Information Systems Laboratory



Reconcile project:
http://reconcile.pl/

Transfer evaluation to semantically similar statements
(claims)

- Human evaluation is at the origin of every automated credibility evaluation task

- Same is true for any semantic or pragmatic task (e.g. translation, image labeling etc.)

- The Big Question: where is the ground truth?

- The answer: ask the crowd or experts

## Supervised Learning

# CASE STUDY 2: DATA INTEGRATION

# Example: Schema Matching

- Integration of heterogeneous data sources
  - Every project on Big Data analysis first has to integrate data from different, <u>heterogeneous</u> data sources
  - One of the long-standing open problems in data management (both industry and research)
- How to find good "matches"?
- How to choose the "best matches"?

**Distributed Information Systems Laboratory**

- **Manual matching**
  - still common practice today

- **Schema matching tools**
  - Based on structural and content features
    - names, domains, structure, values, ...
  - Establish correspondences and rank according to quality
    - Errors are frequent and unavoidable
    - Works well for small schemas

Distributed Information Systems Laboratory

Data integration networks:
different experts may contribute partial matches

Which one would you choose?



Instead of considering only one mapping, consider whole networks of mappings: **leverage knowledge from the network!**

- By combining different matches in a network we can construct evidence for the correctness of those matches

  - For example, a matching contributing to a "bad cycle" less likely to be correct

- Idea: combine all this evidence and use probabilistic reasoning to select the most likely matchings

**variable $x$ to local factor $f$:**

$$\mu_{x \to f}(x) = \prod_{h \in n(x) \setminus \{f\}} \mu_{h \to x}(x)$$

**local factor $f$ to variable $x$**

$$\mu_{f \to x}(x) = \sum_{\sim \{x\}} \left( f(X) \prod_{y \in n(f) \setminus \{x\}} \mu_{y \to f}(y) \right)$$



P. Cudré-Mauroux, K. Aberer and A. Feher. Probabilistic Message Passing in Peer Data Management Systems, ICDE 2006.

- Probabilistic reasoning results in reasonable improvement of matching quality, but
  - a posteriori analysis can only identify potentially bad choices by experts, but not correct them
- Better approach
  - Let experts make better local decisions by providing them information on global consistency and asking targeted questions



Q. V. H. Nguyen, T. Nguyen Thanh, Z. Miklos, K. Aberer and A. Gal et al. Pay-as-you-go Reconciliation in Schema Matching Networks. ICDE 2014.

26

# Minimal Effort User Feedback?

- Asking the right questions is important



Two possible solutions:
{c1,c2,c3} and {c1,c4,c5}

- Ask $c_1$ first
  - → the network is unchanged
  - → no uncertainty reduction.
- Ask $c_2$ first
  - → only 1 solution left
  - → the network becomes certain.

- Idea: optimize information gain with each question

27

Distributed Information Systems Laboratory

- Information gain ordering strategy achieves savings of up to 48% user effort compared to random ordering

- Outperforms the baseline with an average difference of 15% (precision) and 14% (recall)

Distributed Information Systems Laboratory

- Data Integration is a task that combines human and machine intelligence

- The Big Question: How to minimize human effort and maximize information gain?

## **Active Learning**

# CASE STUDY 3: SOCIAL MEDIA ANALYSIS

Distributed Information Systems Laboratory

- Social Media (e.g. Twitter) contains many (hidden) signals on the public perception of issues of general interest
  - nutrition, health, politics, environment etc.
- Goal: identify influencers, their communities, their topics of interest and their stance towards given issues
- Methods
  - Semantic content analysis to capture and classify relevant content
  - Social network analysis to capture and analyze social influence

1. Describe the interest (keywords, users, time, geographic)
2. Select (or collect) the data
3. Extract the key Concepts, Entities and Categories
4. Identify Topics and Communities
5. Select relevant Issues, Influencers and Events
6. Produce insights (correlations)

Distributed Information Systems Laboratory

EPFL
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Distributed Information Systems Laboratory

*Input:*
~50 Mio tweets

*Analysis method:*
Word embedding
(word2vec)

**Panel 1 — food**

Exploring the topic "food" the system suggest related terms. We find terms related to food ingredients

Choose action

{diet, {diet}, 41 583}
{chemicals, {chemicals}, 8725}
{cheese, food, 11 969}
{fish, food, 6832}
{fat, {fat}, 32 555}
{gluten free, {gluten free}, 4959}
{calories, {calories}, 8865}
{candy, {candy}, 4625}
{grain, food, 2961}
{bread, {bread}, 11 656}
{cereal, {cereal}, 3196}
{grains, {grains}, 1861}
{cereals, {cereals}, 2730}
{cooking, {cooking}, 9299}
{feed, {feed}, 15 741}
{diets, {diets}, 6333}
{fruits, {fruits}, 5832}
{ground beef, {ground beef}, 300}
{eggs, {eggs}, 3554}
{drink, {drink}, 24 664}
{cspi, {cspi}, 471}
{big organic,
{butter, {butt
{apples, {ap
{conventiona

**Discovery of interesting dimensions**

**Panel 2 — food ingredients / food / nutrition**

We create a category "food ingredients". The system proposes more related terms.

Choose action

{fat, food ingredients, 32 555}
{calories, food ingredients, 8865}
{sugar, {sugar}, 34 228}
{fructose, {fructose}, 3915}
{sugars, {sugars}, 1484}
{sodium, {sodium}, 2283}
{saturated fat, {saturated fat}, 1146}
{fiber, {fiber}, 2222}
{fats, {fats}, 3886}
{cals, {cals}, 437}
{salt, {salt}, 5759}
{grains, food, 1861}
{nutrients, {nutrients}, 2517}
{chemicals, food ingredients, 8725}
{saturated fats, {saturated fats}, 262}
{low fat, {low fat}, 822}
{added sugar, {added sugar}, 590}
{protein, {protein}, 8042}
{processed foods,
  {processed foods}, 1019}
{trans fats, {trans fats}, 680}
{calorie, {calorie}, 4801}

**Panel 3 — food ingredients / food / nutrition**

We select all terms that are related. We may repeat this step.

Choose action

{low fat, {low fat}, 822}
{added sugar, {added sugar}, 590}
{protein, {protein}, 8042}
{processed foods,
  {processed foods}, 1019}
{trans fats, {trans fats}, 680}
{hfcs, {hfcs}, 1172}
{calcium, {calcium}, 2453}
{trans fat, {trans fat}, 732}
{carb, {carb}, 2232}
{calorie, {calorie}, 4801}
{added sugars, {added sugars}, 233}
{toxins, {toxins}, 2033}
{sweeteners, {sweeteners}, 2030}
{artificial sweeteners,
  {artificial sweeteners}, 978}
{sat fat, {sat fat}, 172}
{additives, food ingredients, 1750}
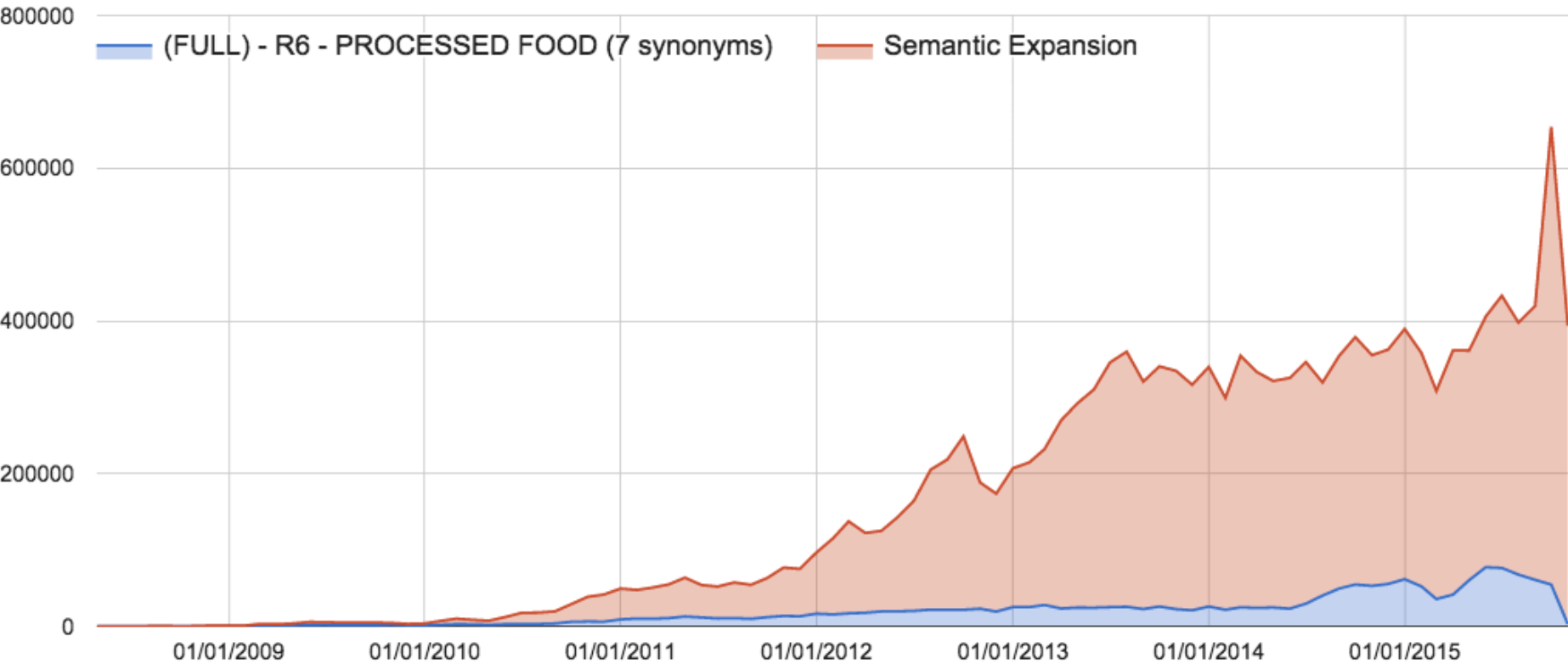{antioxidants, {antioxidants}, 998}
{caffeine, {caffeine}, 802}
{vitamins, {vitamins}, 2312}
{pesticides, {
{diet soda, {
{fruit juice, {f
{bacteria, {ba
{hormones, {hormones}, 1262}

**Instantiation of dimensions with terminology**

**Panel 4 — food / food ingredients / nutrition**

Finally we have a clean list of all terms on food ingredients.

Choose action

{sugar, food ingredients, 34 228}
{fat, food ingredients, 32 555}
{sodium, food ingredients, 2283}
{fructose, food ingredients, 3915}
{fiber, food ingredients, 2222}
{added sugar, food ingredients, 590}
{protein, food ingredients, 8042}
{salt, food ingredients, 5759}
{nutrients, food ingredients, 2517}
{saturated fat, food ingredients, 1146}
{calcium, food ingredients, 2453}
{chemicals, food ingredients, 8725}
{hfcs, food ingredients, 1172}
{trans fats, food ingredients, 680}
{caffeine, food ingredients, 802}
{artificial sweeteners,
  food ingredients, 978}
{sweeteners, food ingredients, 2030}
{carb, food ingredients, 2232}
{trans fat, food ingredients, 732}
{12}
{s, 998}
{1}
{additives, food ingredients, 1750}

Distributed Information Systems Laboratory



Processed Food - 7 Synonyms vs. Semantic Expansion

Using a semantically expanded terminology increases coverage significantly!

Distributed Information Systems Laboratory

- The system clusters the terms on food ingredients according to similarity

- The expert sees
  - A clear distinction between positive and negative terms
  - Distinction between natural and artificial ingredients
  - Clusters of related terms, e.g. vitamins, additives etc.

- We may use this to create sub-categories of interest

![EPFL - ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE]

Distributed Information Systems Laboratory

# Analyzing social interactions we can identify clear communities



Level 0

Collection from cluster | Construct topic | Remove from collection | Filter...
Minimal size: 0 | Maximal size: 5000000

Click to select cluster and see it's documents and concepts, double-click to open subclusters

Info | Users(28101) | Tweets(81720) | Documents(7274) | Concepts(100) | Hashtags(100)
Keywords(205)

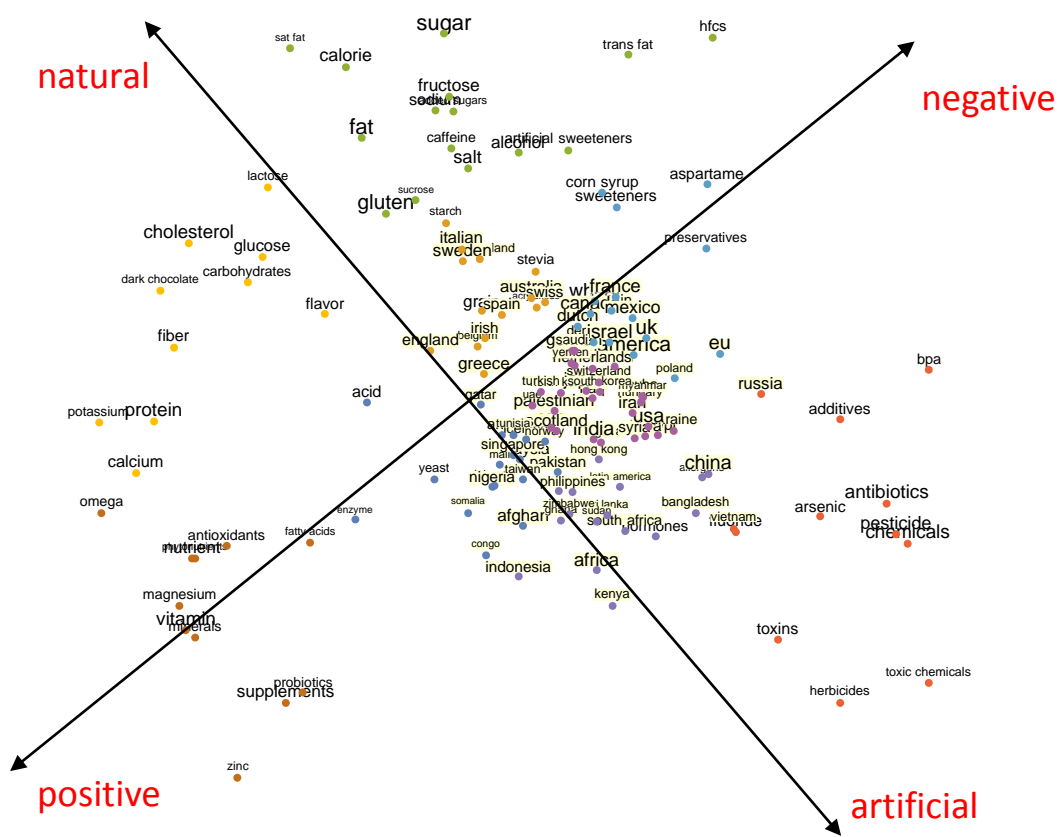| # | Name | Description | Tweets | Retweets | Influence (RT+RE+Mentions) |
|---|------|-------------|--------|----------|---------------------------|
| 1 | 8extremes | Sharing Peace & Understanding with Compassion & Gratitude. Having Love & Respect for all living things. Susan Elaine Los Angeles CA | 872 | 214 | 2367 |
| 2 | GMOFreeUSA | We are a national group, educating consumers about the potential hazards of Genetically Engineered foods. | 78 | 25 | 1679 |
| 3 | OrganicLiveFood | We provide info regarding #organic/sustainable foods, nutritional facts 4improving #health, #raw foods, #superfoods, #herbs, #GMO harms, #bees & #pesticides | 606 | 28 | 1439 |
| 4 | RachelsNews | KIDS RIGHT TO KNOW(Founder)16 DEBATED Shark/Dragon https://t.co/pn1AzZnHKI TEDxTO 2014-SPEAKER https://t.co/3OkImzLnch Huffington Blogger, Teen Earth Activist | 1347 | 223 | 1412 |
| 5 | kevinfolta | Land-grant scientist exploring ways to make better food with less input, also learning and teaching how to effective | 394 | 11 | 1363 |
| 6 | TheGOPJesus | Politicians special pla maybe not | | | |
| 7 | GMWatch | Countering | | | |
| 8 | MonsantoCo | Monsanto others to a sustain the | | | |
| 9 | geneticmaize | Mom of Ad https://t.co sustainable ag. Former US Army public health. Ecomodernist. Words are mine. | | | |
| 10 | SSF_BERF_DEFM | SSF = Support Small Farms -- BERF = By Eating Real Food -- DEFM = Don't Eat Factory Meat | 914 | 0 | 888 |

For each community we can identify
- Their influencers
- Their main concepts
- Potentially new interesting terminology

37

Distributed Information Systems Laboratory

- Machine learning applied to Big Data can reveal surprising hidden structures with valuable insights

- Big questions:
  – How to guide the machines to the right data and analysis
  – How to make the resulting structures human-interpretable

# **Unsupervised Learning**

Distributed Information Systems Laboratory

- Big Data has impressive potential to create insights and solve hard problems

- Human intervention in the analysis processes is essential for obtaining meaningful results

- Three main types of intervention

  - A priori: supervised learning

  - Interactive: active learning

  - A posteriori: unsupervised learning

- No one size fits all: their specific implementation depends strongly on the use case