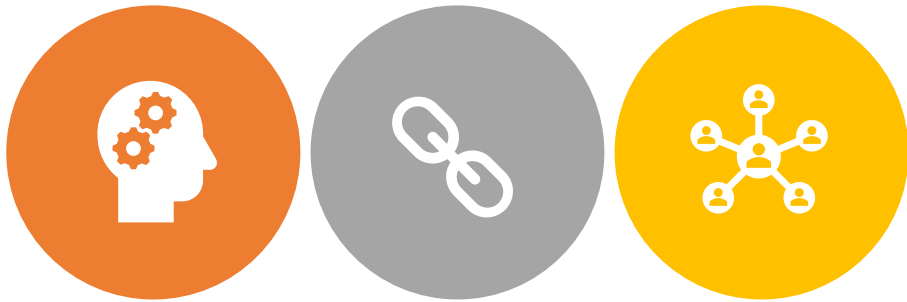# Semantic drift of cited references in the medical literature

Gustaf Nelhans & Johan Eklund

Swedish School of Library and Information Science, University of Borås

24th Nordic Workshop on Bibliometrics and Research Policy, 2019-11-28

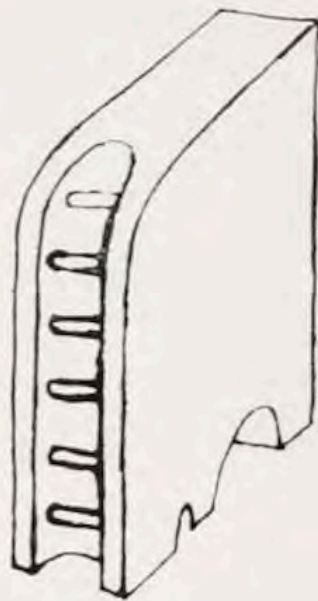# Combining machine learning with citation (linked) data

Theoretical conceptualisation:

"Can we infer the 'knowledge potential' in cited references by identifying/evaluating the citing context around a reference?"

*Hjørland, B. (1992):* The concept of 'subject' in Information Science, *Journal of Documentation 48*(2), 172-200

ONTINUUM" is composed of a
 forms. In full scale, the artist intends
stone, concrete or fiberglass) in a play-
n can clamber over them. In some cases
 inside (see back page) where rungs
ced to help them. Although these ab-
etically pleasing in themselves, the most
al thing about them is the way in which
d and arranged them so that when they
ny one of the apertures in the sculptures,
 forms are grouped together to form a
: a giraffe, a man standing on his head,
 forth, depending on the viewpoint se-
ugh the aperture not only focuses one's
hat make an image but at the same time
eces from view. The playground group
 delightful experience for children (and
lso provide an opportunity for them to
pture.

ontinuum" the artist has had to solve
ems. The pieces forming the composite
ced in an exact line of projection origi-
 aperture. In determining their relative
m the hole and from each other had to
t as well as the height of the viewpoint.
 piece containing the perforation must
ther composite figure when seen from
e perforation had to be placed so that
rt of the composite image or not inter-
ay. In speaking of the complexity of his
 says, "The general concept of the group
oice, and in the early stages of design

The above drawings indicate
sculptures so that children c

an element of free
lationship of the
working with m
those governi
in this disc
musical
and i
w

# Toward a machine-based understanding of text…

...combined with citation data

# Research programme: 'References as words'

- We propose an approach for representing documents by their ***cited references*** rather than by their words

- Scientific ***references*** share the following qualities with symbols (e.g. references) in language:
    - they are ***reusable*** units of meaning
    - they import ***meaning*** into the ***context*** in which they are used

C.f. Small 1978: 'Citations as symbols' SSS.

Topics

Documents

Topic proportions & assignments

Topic modeling

Bag of words

# Topical analysis of reference contexts

In the generated topic model, each word is associated with a probability distribution of topics

congue risus feugiat **ref264** tincidunt lorem nullam

For each reference, a symmetric context window of size $k$ is used as a pseudo-document, and the most probable topic is calculated for that context window

congue risus feugiat **ref264** tincidunt lorem nullam

Eklund, J., & Nelhans, G. (2017). Topic modelling approaches to aggregated citation data. Presented at the *22nd International Conference on Science and Technology Indicators (STI), Paris*, September 6-8, 2017.

**PubMed.gov**
US National Library of Medicine
National Institutes of Health

PubMed

Advanced

Format: Abstract ▾

Send to ▾

**Full text links**

Cochrane Library

# Physical training for asthma.

Carson KV, Chandratilleke MG, Picot J, Brinn MP, Esterman AJ, Smith BJ.

**Save Items**

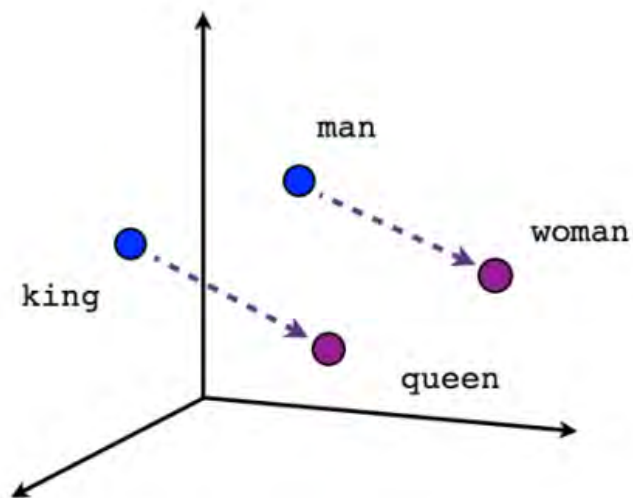## Description of the condition

Asthma, a chronic respiratory condition affecting 300 million people globally (Masoli 2004), causes inflammation of the lungs as well as structural and functional remodelling of the airways. It is characterised by recurrent attacks of breathlessness and wheezing with varying degrees of frequency and severity, which is caused by swelling of the bronchial tubes resulting in airflow limitation (WHO 2011). Although the causes of asthma are not completely understood, risk factors are known to include inhaling asthma triggers such as allergens, tobacco smoke and chemical irritants. Asthma is incurable and the prevalence is increasing, particularly in children and young adults (Pawankar 2012), however appropriate management can control the disorder and enable people to enjoy a high quality of life (WHO 2011).

Asthma, a chronic respiratory condition affecting 300 million people globally ( aref15080825 ), causes inflammation of the lungs as well as structural and functional remodelling of the airways. It is characterised by recurrent attacks of breathlessness and wheezing with varying degrees of frequency and severity, which is caused by swelling of the bronchial tubes resulting in airflow limitation (WHO 2011). Although the causes of asthma are not completely understood, risk factors are known to include
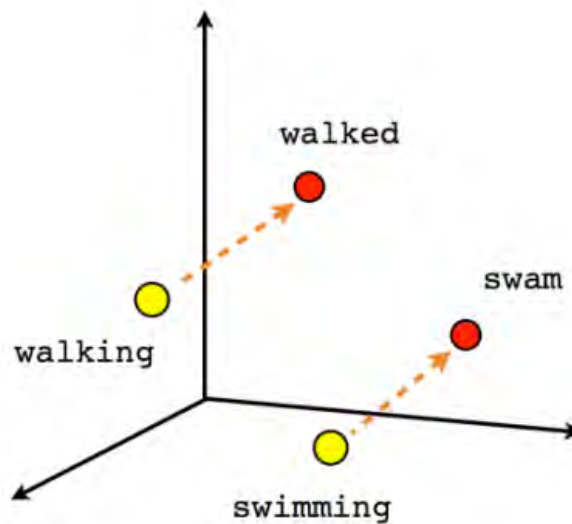
*asthma a chronic respiratory condition affecting million people globally aref causes inflammation of the lungs as well as structural and functional remodelling of the airways*

**Topic 346 (0.8149):** asthma, copd, allergic, airway, disease, fev, ige, respiratory, lung, symptoms
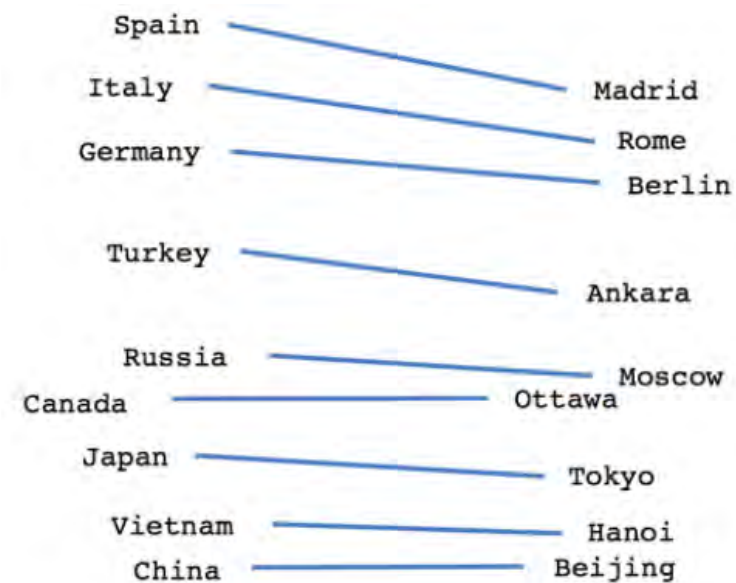
**Topic 78 (0.0689):** pressure, lung, pulmonary, respiratory, gas, lungs, ventilation, volume, breathing, alveolar

Male-Female      Verb tense      Country-Capital

# Word embedding

# Exploring the use of references as words

**RQ1.** Can we identify the degree of *topical heterogeneity* of a subset of investigated cited references?

**RQ2.** Can we *infer the presence of a cited reference* in a given text using our trained model? Correspondingly: *can we reconstruct the context of a reference* in a text?

**RQ3.** Can we identify the *semantic drift* in cited references over time?

# Data

**2.3M Full-text open access articles in the Europe PMC data set:**

*Extract* **text within <BODY> </BODY> in XML files**

**For every reference matched in PubMed:**

Replace in-text reference with PMID.

*Cleaning*: **tokenization, remove punctuation, numbers (except PMID:s), stopwords**

**Extract ~7M reference contexts "citances" around each reference**

20 words before/after the reference: in the following form:**[20W] PMID [20W]**

**Additional analysis steps:**

Identify *publication type* assigned by PubMe

Add *publication year* to each citance.

# Citances: [20W] PMID [20W]

3549392  16221197       accumulating epidemiologic evidence suggests hypovitaminosis associated increased risk cardiovascular events experimental data generally support hypothesis vitamin protective role cardiovascular health **16221197** paper examine relevance omega vitamin cardiology provide update clinical trial results dietary sources pufa fish major food source long chain       [2012     article type]

3549392  16221197       **...**accumulating epidemiologic evidence suggests hypovitaminosis **is** associated **with** increased risk **of** cardiovascular events**. [E]**experimental data generally support **the** hypothesis **that** vitamin **[X plays a]** protective role **in** cardiovascular health **16221197** paper examine relevance omega-**3...** vitamin cardiology provide update clinical trial results dietary sources pufa fish major food source long chain       [2012     article type]

# Results

RQ1: topical heterogeneity

**RQ2. Can we *infer the presence of a cited reference* in a given text using our trained model? Correspondingly: can we *reconstruct the context of a reference* in a text?**

# Bayesian.py: [Terms: gender, differences, smoking]

| PMID | Title |
|------|-------|
| 12627467 | **Gender**-specific molecular heterosis and association studies: dopamine D2 receptor gene and **smoking**. |
| 15831689 | Trends in smoking behaviour between 1985 and 2000 in nine European countries by education. |
| 17220338 | Global DNA methylation level in whole blood as a biomarker in head and neck squamous cell carcinoma. |
| 12186222 | Cigarette **smoking** among adults--United States, 2000. |
| 8643986 | **Gender differences** in health: are things really as simple as they seem? |
| 12200093 | Trends in smoking, diet, physical exercise, and attitudes toward health in European university students from 13 countries, 1990-2000. |
| 9150319 | **Gender difference** in **smoking** effects on lung function and risk of hospitalization for COPD: results from a Danish longitudinal population study. |
| 15081207 | **Gender differences** in health: a Canadian study of the psychosocial, structural and behavioural determinants of health. |
| 12241535 | *Determination of age-related increases in large artery stiffness by digital pulse contour analysis.* |
| 12573363 | *Paraoxonase and coronary heart disease.* |

# 'Pointwise mutual information'

- [Given PMID: 12627467]
  - "Gender-specific molecular **heterosis** and association studies: dopamine D2 receptor gene and smoking."

| [PMID: 12627467] | |
| --- | ---: |
| **heterosis** | 0.798637 |
| heterozygotes | 0.64744 |
| alarmingly | 0.60751 |
| negatives | 0.606381 |
| homozygotes | 0.603385 |
| directions | 0.552902 |
| discriminative | 0.547734 |
| inconsistencies | 0.533126 |
| drd | 0.500498 |
| taq | 0.46735 |

# RQ2 reference embedding word2vec

# Reference embedding

*"Subjects in themselves must thus be defined as the epistemological potentials of documents" (Hjørland, 1992)*

# Thank you!