Bibliothèque
universitaire de médecine

# Mine and combine

Text Mining Tools Used for Search Term
Identification
Meet and Greet

UNIL | Université de Lausanne
Faculté de biologie et de médecine

Cécile Jaques
Jolanda Elmers

# Overview

- Brief introduction and presentation

  Focus on     PubMed PubReMiner

                    TerMine

                    Yale MeSH Analyzer

- Hands-on session
- Discussion together
- More tips and tricks from our experience

UNIL | Université de Lausanne
CHUV
Faculté de biologie et de médecine

# Questions

Who has been involved in developing a search strategy for a systematic review?

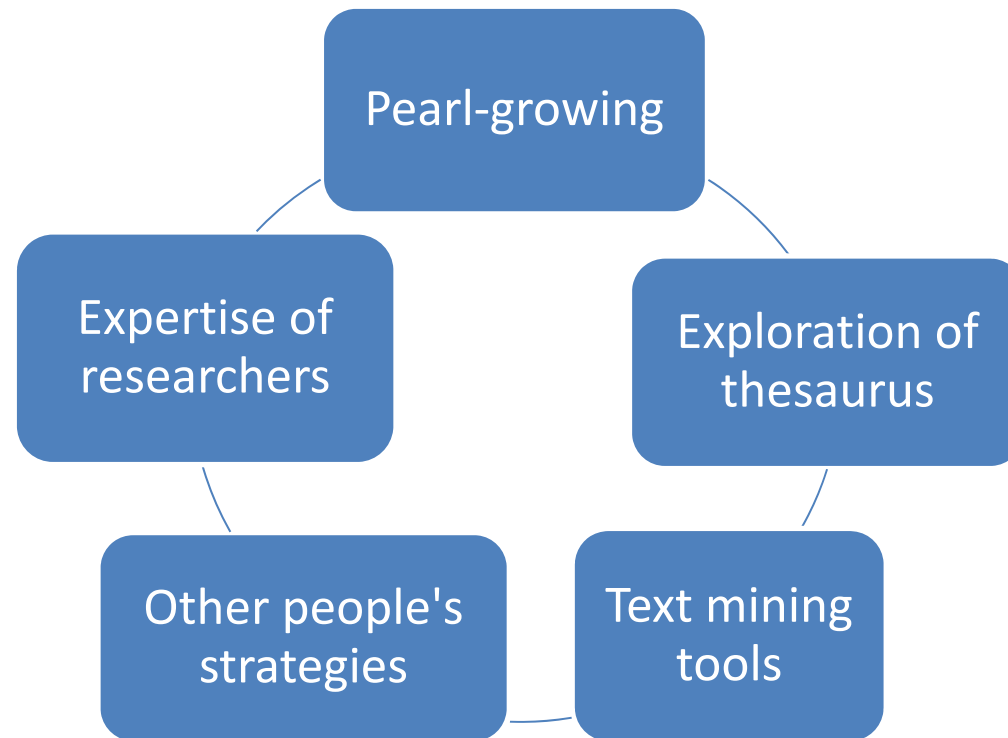Who has already tested the following tools (or one of them)?
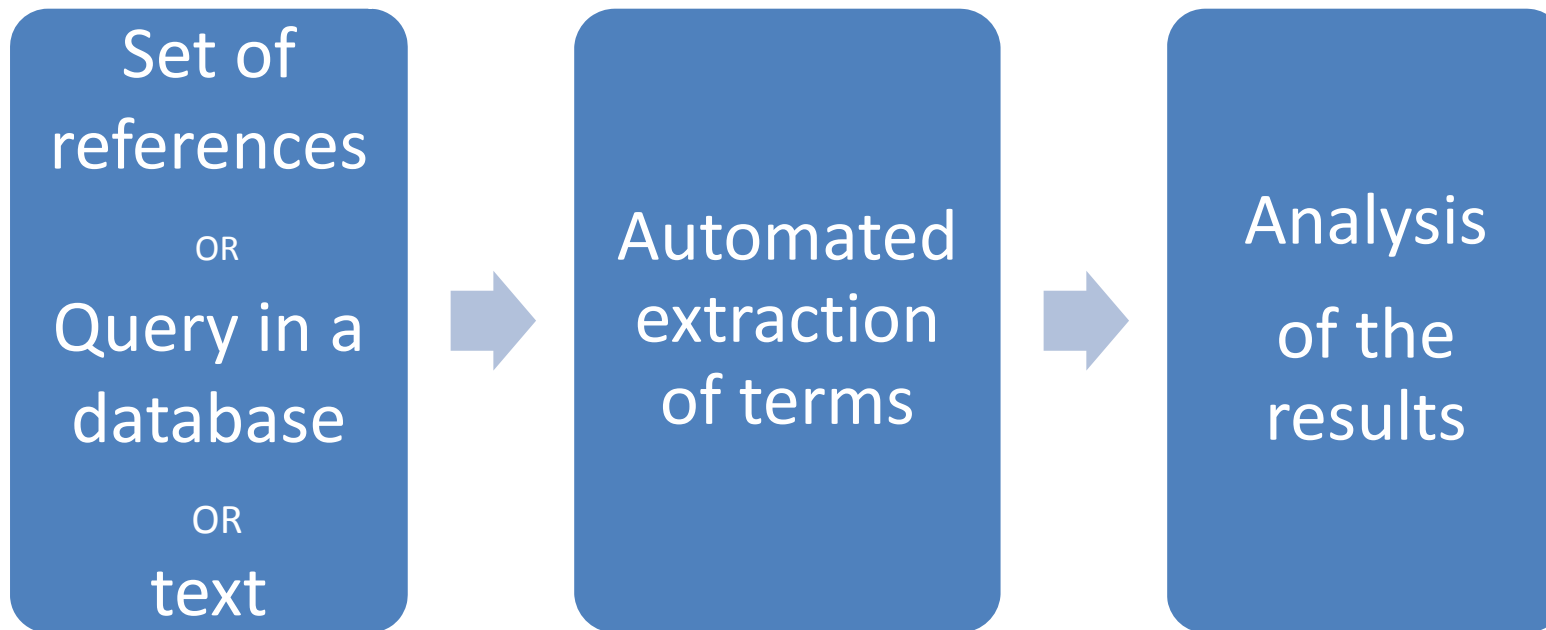>    PubMed PubReMiner
>    TerMine
>    Yale MeSH Analyzer

Who uses these tools (or one of them) on a regular basis?

# Search term identification methods

# Text mining tools

| Set of references | | Automated extraction of terms | | Analysis of the results |
|---|---|---|---|---|
| OR | | | | |
| Query in a database | → | | → | |
| OR | | | | |
| text | | | | |

Different programmes offer different levels of analysis

UNIL | Université de Lausanne
CHUV
Faculté de biologie et de médecine

# Overview of the tools

| PubMed PubReminer | Frequency count of words in a set of articles, based on a query in PubMed |
|---|---|
| **Yale MeSH Analyzer** | MeSH analysis grid to identify MeSH Terms from a set of PubMed articles |
| **Termine** | Word combinations are extracted from raw text from any database, based on statistical and linguistic analysis |

UNIL | Université de Lausanne
Faculté de biologie et de médecine

# PubMed PubReminer

- Uses frequency tables for terms to display results of search queries in PubMed

- Offers a quick view of the words which are occurring most frequently in your results set

# PubMed PubReminer



Detailed analysis of PubMed Search results

**Enter your PubMed Query**

Start remining PubMed for:
```
("Air Pollution"[Mesh] OR air pollution[tiab]) AND
("Cardiovascular Diseases/epidemiology"[Mesh] OR
(cardiovascular disease*[tiab] AND morbidity))
```

Fieldtype: All

Publicationtype: All

FromDate: YYYY/MM/DD (Optional)

ToDate: YYYY/MM/DD (Optional)

AbstractLimit: 10000

Start PubReMiner    Reset

**Lookup a human gene and use all its synonyms**

Lookup Gene:

Search Gene    Reset

Any query that can be processed by PubMed

Up to 10 000 abstracts

Click on start to proceed to the analysis

UNIL | Université de Lausanne    CHUV    Faculté de biologie et de médecine

# PubMed PubReminer - results

# PubMed PubReminer - results

# refers to the number of references the term appears in

Count refers to the number of times the term appears in total

| # | Count | OR | Word |
|---|---|---|---|
| 1954 | 5811 | ☐ | POLLUT * |
| 1706 | 5830 | ☐ | EFFECT * |
| 1687 | 1880 | ☐ | HUMAN * |
| 1661 | 5082 | ☐ | DISEASE * |
| 1624 | 6940 | ☐ | AIR |
| 1408 | 3099 | ☐ | EPIDEMIOLOGY |
| 1395 | 3617 | ☐ | CARDIOVASCULAR |
| 1316 | 4787 | ☐ | EXPOSURE * |
| 1300 | 2708 | ☐ | ADVERSE * |
| 1297 | 3904 | ☐ | RISK * |
| 1230 | 3190 | ☐ | INCREASE * |
| 1119 | 2281 | ☐ | DATA |
| 1106 | 2223 | ☐ | ASSOCIATE * |
| 1073 | 2748 | ☐ | ASSOCIE * |
| 1068 | 2848 | ☐ | ANALYSE * |
| 1058 | 5175 | ☐ | MORTALITY |

| # | OR | Mesh | |
|---|---|---|---|
| 2954 | - | / epidemiology | P |
| 2378 | - | / adverse effects | P |
| 1652 | ☐ | Humans | P |
| 1380 | ☐ | / analysis | P |
| 1240 | ☐ | / mortality | P |
| 1110 | ☐ | / etiology | P |
| 923 | ☐ | / statistics & numerical data | P |
| 844 | ☐ | Female | P |
| 826 | ☐ | Male | P |
| 807 | ☐ | Air Pollution | P |
| 763 | ☐ | Cardiovascular Diseases | P |
| 650 | ☐ | Middle Aged | P |
| 634 | ☐ | Aged | P |
| 614 | ☐ | Air Pollutants | P |
| 570 | ☐ | / toxicity | P |
| 555 | ☐ | Air Pollution/adverse effects | P |

P will process a search in PubMed for the chosen term

A click on the term will add it to the query in PubReminer

10

UNIL | Université de Lausanne
CHUV
Faculté de biologie et de médecine

# PubMed PubReminer
# advanced search options

Click on a hyperlink to add that element to your query and Re-Mine or select terms (OR boxes) and press 'Search Again'
Click on the **P** to directly goto PubMed and view ALL references for that element
Save the results as a txt-file

**Operator:** AND ▾ **Merge similar words:** YES ▾ **Minimalcount:** 2 Search Again

Add selected items from the results analysis to the query
available operators : AND / NOT

does not work with MeSH terms with multiple subheadings

Only for the word column

# TerMine

- Tool used for term recognition
  - recognises automatically candidate multiword terms from documents.

- It annotates raw text from any database

- It recognises acronyms

# TerMine Web Demonstration

**Web Demonstration**

○ Plain text (Only ASCII characters allowed)

Different options
of introducing
text
(2MB maximum)

○ Local text file (*.txt file in ASCII encoding or *.pdf file; 2MB maximum)
[Choisissez un fichier] Aucun fichier choisi

○ URL (HTML or PDF content; 2MB maximum)

POS tagger: [GENIA Tagger version 2.1 ▾] ☑ Preserve break lines

[Analyze] [Clear] [Try (NaCTeM sample)] [Try (MEDLINE sample)]

Check that the
selected tagger is
GENIA tagger.
This tagger is
specifically tuned
for biomedical text
such as MEDLINE
abstracts

13

UNIL | Université de Lausanne    CHUV
Faculté de biologie et de médecine

# TerMine results

Choose the in table presentation

## TerMine (C-value) analysis

Service

Found **1526** terms in 12.6 seconds - all terms (in table) (in text) - threshold: [ 0 ] Apply

. 1 Orthopedics. 2017 Mar 1 ; 40 ( 2 ) : e242-e247. doi : 10.3928/01477447-20160901-03 Epub. 2016 Sep 9. . Radial Shaft Reconstruction With an Intercalary Endoprosthesis Following. Resection of Metastatic Tumor. . Gibson PD , Ippolito JA , Benevenia J. . Improvements in imaging and treatment of musculoskeletal tumors have increased. the variety of options for reconstruction following joint-sparing diaphyseal. resection The purpose of this case series was to show that reconstruction of. malignant tumors of the radial shaft with an intercalary prosthesis may be an. option for patients with segmental bone loss Three consecutive patients. underwent wide resection of the radial diaphysis followed by reconstruction with. a custom intercalary prosthesis A custom intercalary prosthesis with lap joint. design was used in all 3 cases Mean follow-up was 18 months ( range , 9-25. months ) All patients were weight bearing as tolerated 1 week postoperatively At. the most recent follow-up , patients ' mean elbow flexion and extension arc was. 137 ? ? ( range , 130 ? ? -140 ? ? ) At the forearm , mean supination was 60 ? ? ( range , 30 ? ? -90 ? ? ). and mean pronation was 70 ? ? ( range , 60 ? ? -90 ? ? ) At the wrist , mean palmar flexion. was 80 ? ? ( range , 70 ? ? -90 ? ? ) and mean dorsiflexion was 80 ? ? ( range , 70 ? ? -90 ? ? ) All. patients reported minimal to no pain and no significant functional limitations. Mean Musculoskeletal Tumor Society score was 26/30 ( 87 % ) Reconstruction with an. intercalary prosthesis is a viable option for patients with metastatic disease of. the radial shaft All patients had satisfactory results and early return to. function ; none required return to the operating room Possible advantages of. reconstruction with an intercalary prosthesis compared with reconstruction with a. bone graft or polymethylmethacrylate osteosynthesis include early return to. function and minimal weight-bearing restrictions postoperatively. [ Orthopedics. 2017 ; 40

14

# TerMine

## Results in table format

- Score indicates the c-value. This value is calculated using automatic term recognition, based on the following characteristics:

  - the occurrence frequency of the candidate term
  - the frequency of the candidate term as part of other longer candidate terms
  - the number of these longer candidate terms
  - the length of the candidate term

| Rank | Term | Score |
|------|------|-------|
| 1 | author information | 43.542858 |
| 2 | lateral elbow pain | 40.488586 |
| 3 | elbow pain | 30.97436 |
| 4 | lateral elbow | 27.733334 |
| 5 | upper extremity | 15.678572 |
| 6 | tennis elbow | 13.8 |
| 7 | musculocutaneous nerve | 13.125 |
| 8 | memorial chiropractic college | 13.07594 |
| 9 | pain intensity | 12.909091 |
| 10 | musculoskeletal disorder | 12.428572 |
| 11 | canadian memorial chiropractic | 12.362708 |
| 12 | chiropractic college | 12.25 |
| 13 | canadian memorial | 11.8 |
| 14 | ontario institute | 11 |
| 15 | musculoskeletal pain | 10.894737 |
| 16 | distal bicep | 10.777778 |
| 17 | uoit-cmcc centre | 10 |
| 18 | lateral epicondylitis | 9.692307 |
| 19 | nerve injury | 8.846154 |
| 20 | cutaneous nerve | 8.357142 |
| 21 | research associate | 8 |
| 21 | computer user | 8 |
| 21 | cochrane database syst rev. | 8 |
| 24 | forearm muscle | 7.9375 |
| 25 | clinical research excellence | 7.924812 |

UNIL | Université de Lausanne
Faculté de biologie et de médecine

# Yale MeSH Analyzer

- Creates a grid which displays the ways articles are indexed in Medline

- For each article, MeSH Terms are sorted and grouped alphabetically

- Choice to include author keywords, titles and abstracts in the analysis grid

UNIL | Université de Lausanne

CHUV

Faculté de biologie et de médecine

# Yale MeSH Analyzer



Introduce
list of PMID

Only
20 articles
at a time !

# Yale MeSH Analyzer

| PMID | 27610702 | 27281378 | 27179317 | 26130104 | 25481709 | 25063413 |
|---|---|---|---|---|---|---|
| **Title** | Radial Shaft Reconstruction With an Intercalary Endoprosthesis Following Resection of Metastatic Tumor. | Playing-Related Musculoskeletal Problems Among Professional Orchestra Musicians in Scotland: A Prevalence Study Using a Validated Instrument, the Musculoskeletal Pain Intensity and Interference Questionnaire for Musicians (MPIIQM). | Is synergistic organisation of muscle coordination altered in people with lateral epicondylalgia? A case-control study. | The effectiveness of exercise for the management of musculoskeletal disorders and injuries of the elbow, forearm, wrist, and hand: a systematic review by the Ontario Protocol for Traffic Injury Management (OPTIMa) collaboration. | Mechanistic experimental pain assessment in computer users with and without chronic musculoskeletal pain. | Sonography of the lateral antebrachial cutaneous nerve with magnetic imaging and anatomic correlation |
| **Journal Title** | *Orthopedics* | *Med Probl Perform Art* | *Clin Biomech (Bristol, Avon)* | *J Manipulative Physiol Ther* | *BMC Musculoskelet Disord* | *J Ultrasound Med* |
| **Author (Year)** | Gibson PD (2017) | Berque P (2016) | Heales LJ (2016) | Menta R (2015) | Ge HY (2014) | Chiavaras MM (2014) |
| **MeSH Headings** | Aged<br>Aged, 80 and over | Adult | Adult<br>Analysis of Variance | Accidents, Traffic<br>Adult | Adult | Adult<br>Aged<br>Aged, 80 and over |
| | Bone Neoplasms / physiopathology<br>Bone Neoplasms / secondary<br>Bone Neoplasms / surgery* | | | | | |
| | Carcinoma, Renal Cell / physiopathology<br>Carcinoma, Renal Cell / secondary<br>Carcinoma, Renal Cell / surgery | | Case-Control Studies | Cooperative Behavior | Computers*<br>Cumulative Trauma Disorders / diagnosis*<br>Cumulative Trauma Disorders / epidemiology | |
| | Diaphyses / surgery* | | | Disease Management | | |
| | | | Elbow Joint / physiology<br>Electromyography | Exercise Therapy / methods* | | Elbow / anatomy & histology<br>Elbow / diagnostic imaging*<br>Elbow / innervation* |

18

# Hands-on session

- Work in groups of 2
- Links and text analysis package on
  http://bit.ly/2xGknrl
- 30'

UNIL | Université de Lausanne
CHUV
Faculté de biologie et de médecine

# Discussion

How would you implement the use of a text mining tool today in your own systematic research?

- – What are the advantages/disadvantages?

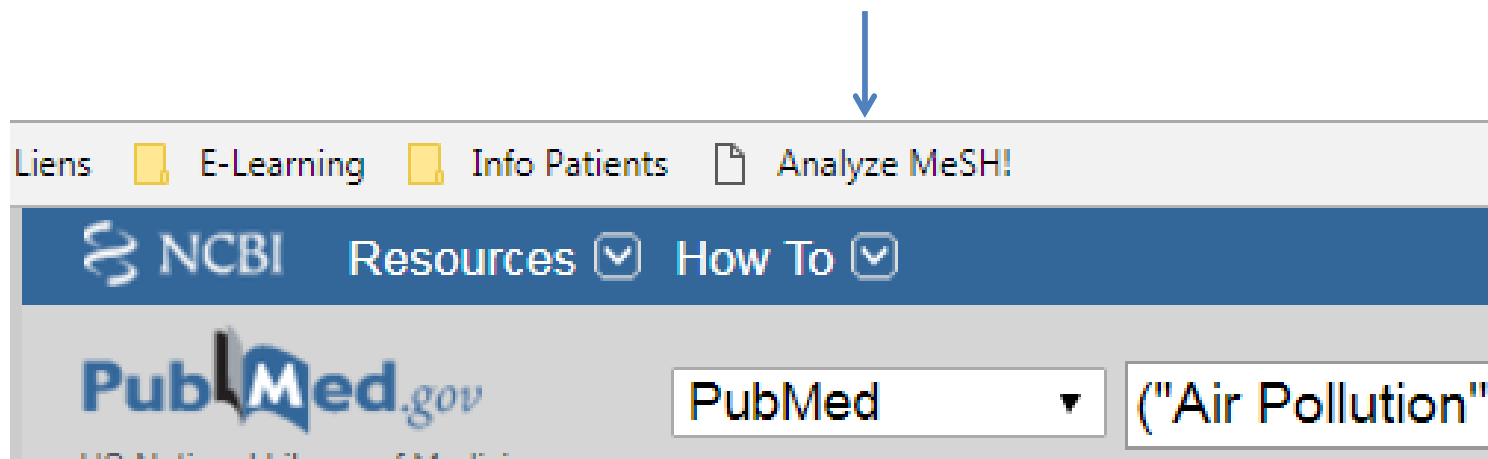- – Which functions do you find useful for your systematic researches?

# Our point of view

| | Useful for | Advantages | Disadvantages |
|---|---|---|---|
| PubReminer | Term identification (ti, ab, MeSH) in gold standard articles provided by the researcher<br>Explore simple PubMed query | Large amount of references<br>Frequency count | |
| Yale Analyzer | MeSH term identification in gold standard articles<br>Starting point for the MeSH analysis with researcher | Easy to use, to read. | Only 20 PMID<br>No frequency count |
| Termine | Phrase identification in references text (ti, ab k eywords)<br>useful when the number of references is quite large | References from other databases | For large amounts of text, necessity to register |

Faculté de biologie et de médecine

# Tips & Tricks - PubReminer

- can help choosing journals for the handsearching process
- Be aware of the Merge function
  - Select No to see all the term variations
  - Select Yes to have a view of the importance of the term root

# Tips & Tricks – Yale analyzer

- can be used directly in PubMed

# Tips & Tricks - Termine

- [Batch Service](): for processing documents larger than 2MB

- create a text file by import / export in Endnote with an export style that keep only relevant text (ti, ab, keywords)

# Report results to researcher

- PubReminer : save as a text file
- Yale MeSH Analyzer : output option - excel
- Termine : create an excel file by copying and pasting the table

UniL | Université de Lausanne

Faculté de biologie et de médecine

# Bibliography

1. Thomas J, McNaught J, Ananiadou S. Applications of text mining within systematic reviews. Res Syn Meth. 2011;2(1):1-14. doi: 10.1002/jrsm.27

2. Belter CW. Citation analysis as a literature search method for systematic reviews. 1 nov 2016;2766-77. doi:10.1002/asi.23605

3. EPC Methods [online]: An Exploration of the Use of Text-Mining Software in Systematic Reviews - NCBI Bookshelf; [Accessed 15 August 2017]. Available from: https://www.ncbi.nlm.nih.gov/books/NBK362044/

4. Tsuruoka Y, Tateishi Y, Kim J-D, Ohta T, McNaught J, Ananiadou S, et al. Developing a Robust Part-of-Speech Tagger for Biomedical Text. In: Advances in Informatics [Internet]. Springer, Berlin, Heidelberg; 2005 [cited 2017 Aug 27]. p. 382–92. (Lecture Notes in Computer Science). Available from: https://link.springer.com/chapter/10.1007/11573036_36

5. O'Mara-Eves A, Thomas J, McNaught J, Miwa M, Ananiadou S. Using text mining for study identification in systematic reviews: a systematic review of current approaches. Systematic reviews. 2015;4(1):5. doi: 10.1186/2046-4053-4-5

UNIL | Université de Lausanne
CHUV
Faculté de biologie et de médecine

Cecile.jaques@chuv.ch
jolanda.elmers@chuv.ch


Bibliothèque universitaire de médecine

http://bium.ch