# Systematic Litterature Review

Christophe Benavent et Olivier Caron

Université Dauphine-PSL - DRM - Acss

February 16, 2023

# Sommaire I

Section 1

Introduction

# Introduction

Systematic reviews are different from traditional literature reviews because they aim to identify all studies that address a specific question. Typically : "How much high are advertising elasticities at the level of the brand.

- so not only a quantitative dimension cause large amount of references.
- Need Text mining methods because we deal mainly with text, focusing on values would be a meta-analyses.
- Toward reproducibility criteria.

In this presentation a pragmatic approach through three questions

- How to constitute the data set ?
- How to analyse Authorship and communities ?
- How to Analyse contents ?

# Narratives versus systematic ?

| Type | Advantage(s) | Disadvantage | Application(s) | Guidelines |
|---|---|---|---|---|
| Systematic Review | 1. Minimized bias<br>2. A-priori protocol<br>3. Defined search and evaluation methods<br>4. Reproducible<br>5. High validity of review conclusions | 1. Must adhere to established guidelines<br>2. Valid literature base required<br>3. Robust (enough) literature to review<br>4. Variation in study methods within reviewed literature may affect results | 1. Identify relevant evidence<br>2. Assess quality of evidence<br>3. Non-biased synthesis of literature<br>4. Interpret evidence in an impartial manner<br>5. Applicable for establishing standards and health policy | PRISMA Guidelines[2] |
| Meta-Analyses - Quantitative | 1. Same as systematic review<br>2. Determine a single estimate of the effect of treatment or management of an illness or event | 1. Data in literature must be homogeneous and available for pooled analysis<br>2. Reliability of literature designs may affect results | 1. Same as systematic review<br>2. Determine best practice for defined topic or event.<br>3. Narrow variations in known data sets. | PRISMA Guidelines[2] |
| Meta-Analyses - Qualitative | 1. Same as systematic review<br>2. Determine major themes or experiences for an event or issue | 1. Variable sampling errors in original literature leads to bias<br>2. Variation in qualitative tools used for original research | 1. Same as systematic review<br>2. Define primary themes and priorities<br>3. Refine future research objectives | PRISMA Guidelines[2] |
| Cochrane Review | 1. Form of systematic review method<br>2. Well defined methodology<br>3. Indexed in the Cochrane Library (open source) | 1. Same as for Systematic Reviews | 1. Same as systematic review<br>2. Determine support for specific treatment<br>3. Determine if evidence exists for defined concept | Cochrane Manual[6] |
| Scoping Review | 1. Use of fluid literature search strategy<br>2. Broader review topics<br>3. May include literature of varied methodologies | 1. Risk of bias due to lack of defined evaluation methods<br>2. Non-specific objectives<br>3. Heterogeneity in literature included | 1. Map available literature in a review field or area<br>2. Literature gap analysis<br>3. Clarification of concept or theory | PRISMA SrR[7] |
| Narrative Review | 1. Researcher determines literature to include<br>2. Less time intensive<br>3. May include literature of varied methodologies<br>4. Interpretive objectives (not structured analysis) | 1. Risk of multiple forms of bias and error<br>2. Unstructured, not reproducible<br>3. May not include all appropriate literature<br>4. Lacks systematic synthesis of literature | 1. Identify theory and frames of thought on a topic<br>2. Summarize a particular study topic<br>3. Justify a research topic | |
| Critical Review | Same as Narrative Review | Same as Narrative Review | 1. Develop perspectives on a topic | |
| Conceptual Review | Same as Narrative Review | Same as Narrative Review | 1. Evaluate general consensus on a topic<br>2. Show gaps of knowledge in literature | |
| State-of-the Art Review | Same as Narrative Review | Same as Narrative Review | 1. Describe current beliefs on a topic | |

Stratton © 2019 Prehospital and Disaster Medicine

Figure 1: Narratives, systematic and others LR Stratton (2019)

# the prisma framework

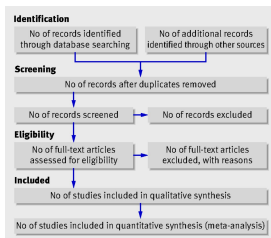A systematic approach to set the article corpus Moher et al. (2009)



Figure 2: Figure 2 : Prisma Process

| Section/Topic | Item # | Checklist Item | Reported on Page # |
|---|---|---|---|
| **TITLE** | | | |
| Title | 1 | Identify the report as a systematic review, meta-analysis, or both. | |
| **ABSTRACT** | | | |
| Structured summary | 2 | Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number. | |
| **INTRODUCTION** | | | |
| Rationale | 3 | Describe the rationale for the review in the context of what is already known. | |
| Objectives | 4 | Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS). | |
| **METHODS** | | | |
| Protocol and registration | 5 | Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number. | |
| Eligibility criteria | 6 | Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale. | |
| Information sources | 7 | Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched. | |
| Search | 8 | Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. | |
| Study selection | 9 | State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis). | |
| Data collection process | 10 | Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators. | |
| Data items | 11 | List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made. | |
| Risk of bias in individual studies | 12 | Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis. | |
| Summary measures | 13 | State the principal summary measures (e.g., risk ratio, difference in means). | |
| Synthesis of results | 14 | Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., $I^2$) for each meta-analysis. | |
| Risk of bias across studies | 15 | Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies). | |
| Additional analyses | 16 | Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified. | |
| **RESULTS** | | | |
| Study selection | 17 | Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram. | |
| Study characteristics | 18 | For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations. | |
| Risk of bias within studies | 19 | Present data on risk of bias of each study and, if available, any outcome-level assessment (see item 12). | |
| Results of individual studies | 20 | For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group and (b) effect estimates and confidence intervals, ideally with a forest plot. | |
| Synthesis of results | 21 | Present results of each meta-analysis done, including confidence intervals and measures of consistency. | |
| Risk of bias across studies | 22 | Present results of any assessment of risk of bias across studies (see Item 15). | |
| Additional analysis | 23 | Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]). | |
| **DISCUSSION** | | | |
| Summary of evidence | 24 | Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., health care providers, users, and policy makers). | |
| Limitations | 25 | Discuss limitations at study and outcome level (e.g., risk of bias), and at review level (e.g., incomplete retrieval of identified research, reporting bias). | |
| Conclusions | 26 | Provide a general interpretation of the results in the context of other evidence, and implications for future research. | |
| **FUNDING** | | | |
| Funding | 27 | Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review. | |

# White and gray

Not only reviewed papers, but also working paper preprints etc.

# Collecting with IA

- Elicit
- Litmaps
- Connected papers

# r environnement

- r + Rstudio + Quarto/beamer to produce this presentation and doing computations.
- You can clone it at github/benaventc.
- Main Packages

```r
library(tidyverse) # the tol Box
library(Rtsne) # 2D magic
library(ggrepel)    #complement de ggplot
library(ggwordcloud)  #complement de ggplot
library(cowplot)

library(udpipe) # annotations
library(quanteda) # un bel ensemble de techniques
library(quanteda.textstats)
library(quanteda.textmodels)
library(quanteda.textplots)
library(fastcluster) #pour aller plus vite
library(ape) #phylo and clustering
library(word2vec) #for embeddings
library(rcrossref)
library(flextable)

theme_set(theme_minimal()+theme(plot.title = element_text(size=12)))
```

Section 2

Data sets acquisition

# Data sets acquisition

- Through citations database, with format (Bib, RIS, json), and API.

# Some Sources

- Google Scholar : harvesting every things
- Crossref : open source
- Scopus : elsevier papers
- Ebsco : business source complete
- Jstor
- NBER
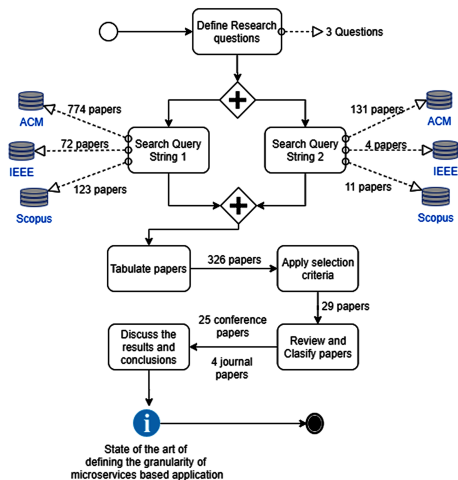- Arxiv and other Psyxiv ou socioxiv, don't forget HAL,
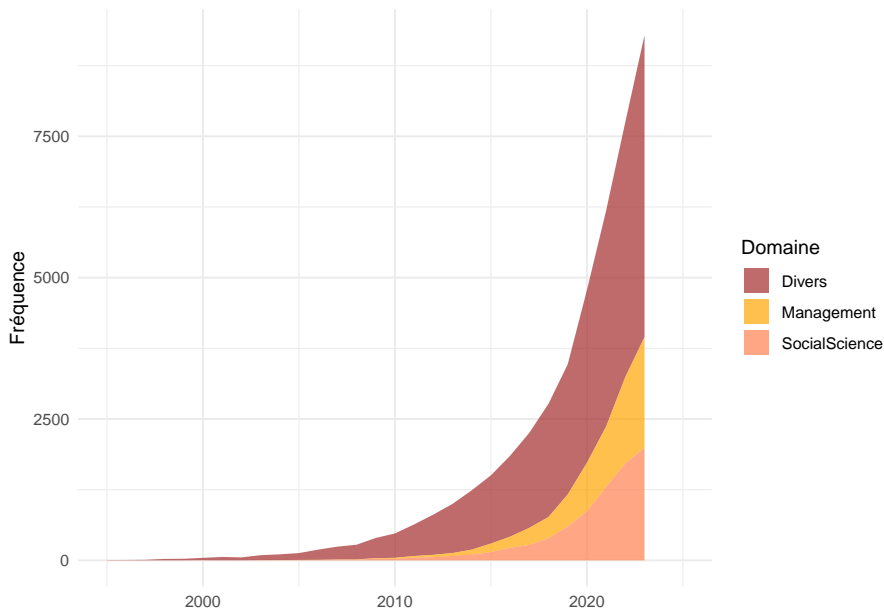
...

# A selection Process



Figure 4: Figure 4 : An SLR process source : https://doi.org/10.7717/peerj-cs.695/fig-1

# A short example with corpus

- Scopus : TITLE-ABS-KEY ("Systematic literature review") -> 37,190 documents
- First health then computing science and a growing concern for social sciences.

# A small Case study : NLP and marketing

A first case : "NLP in Marketing - state of art and evolutions" - Abstract, title, Keywords =("NLP" | "natural language processing" | "Text-Mining"| "text Analysis") & Journal =("Marketing" | "Consumer" )
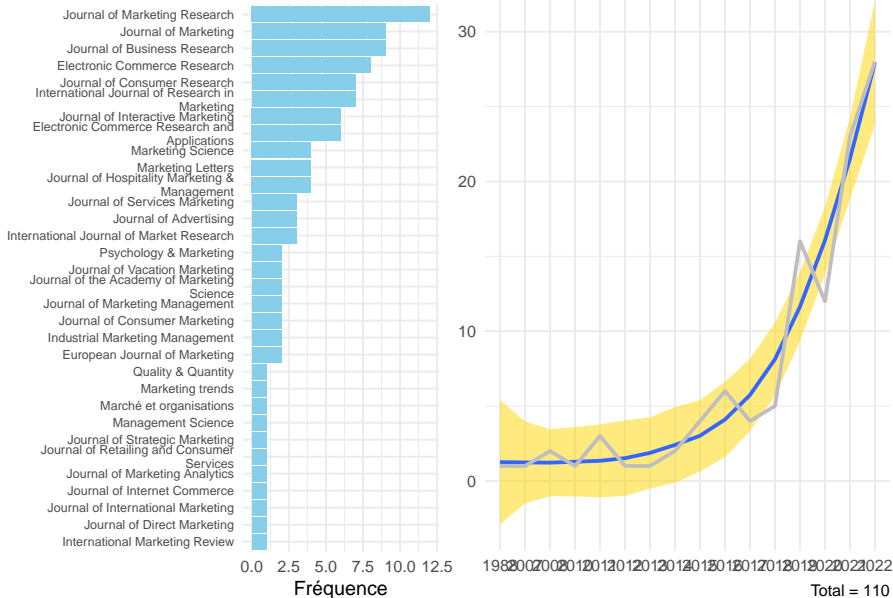
- The references are reported manually, through DOI and maintain in the Zotero collection, then export as datafile.
- Keywords are corrected and completed manually.
- Outcome : after cleaning -> 104 papers.

```
# A tibble: 5 x 9
     id  year auteurs                    title review keywo~1 text  methods fields
  <dbl> <dbl> <chr>                      <chr> <chr>  <chr>   <chr> <chr>   <chr>
1     1  1988 Beard, John D.; William~   Incr~ Journ~ direct~ This~ Readab~ direc~
2     2  2007 Eliashberg, Jehoshua; H~   From~ Manag~ entert~ Movi~ predic~ movie
3     3  2008 Pekar, Viktor; Shiyan Ou   Disc~ Journ~ blogs,~ Auto~ opinio~ hospi~
4     4  2008 Sawyer, Alan G.; Laran,~   The ~ Journ~ readab~ This~ Readab~ publi~
5     5  2010 Nielek, Radoslaw; Wawer~   Spir~ Elect~ Auctio~ An a~ Sentim~ Aucti~
# ... with abbreviated variable name 1: keywords
```
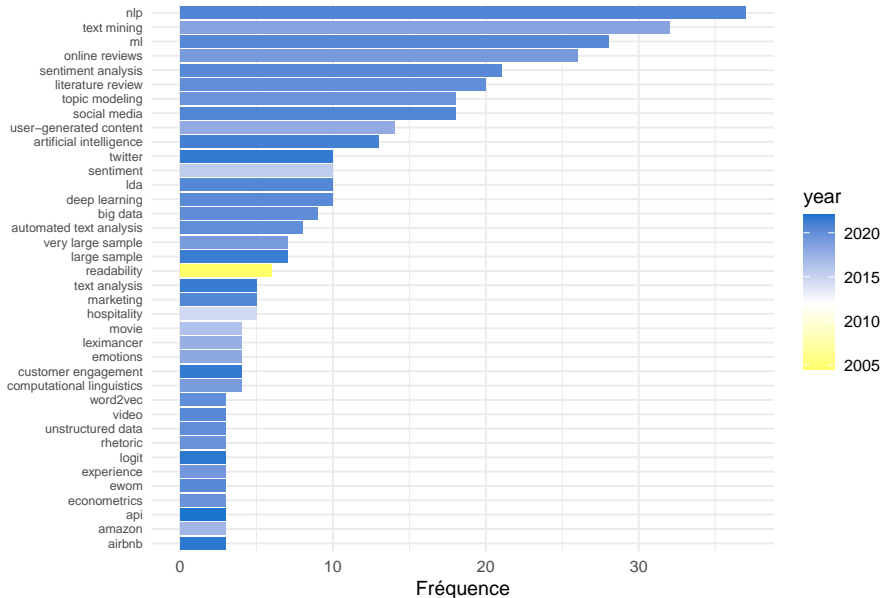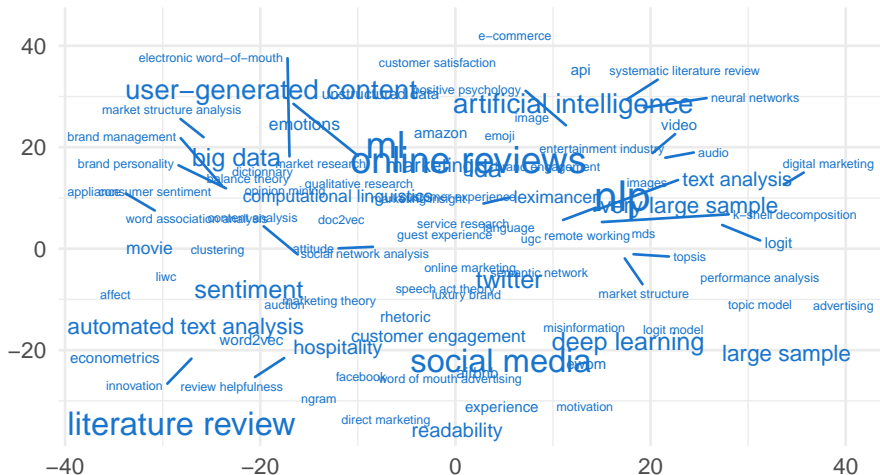
# Marketing et NLP : Nombre d'articles par revue et par an



Journal of Marketing Research
Journal of Marketing
Journal of Business Research
Electronic Commerce Research
Journal of Consumer Research
International Journal of Research in Marketing
Journal of Interactive Marketing
Electronic Commerce Research and Applications
Marketing Science
Marketing Letters
Journal of Hospitality Marketing & Management
Journal of Services Marketing
Journal of Advertising
International Journal of Market Research
Psychology & Marketing
Journal of Vacation Marketing
Journal of the Academy of Marketing Science
Journal of Marketing Management
Journal of Consumer Marketing
Industrial Marketing Management
European Journal of Marketing
Quality & Quantity
Marketing trends
Marché et organisations
Management Science
Journal of Strategic Marketing
Journal of Retailing and Consumer Services
Journal of Marketing Analytics
Journal of Internet Commerce
Journal of International Marketing
Journal of Direct Marketing
International Marketing Review

Fréquence

Total = 110

## Mots clés les plus fréquents

# Projection Tsne des mots clés

# Using API

The best way to operate is to work through API :

- it it prevent from errors
- it's precise
- it support standard formats (bib, ris ...)

Main sources

- Crossref
- Scopus
- [Hal](https://api.archives-ouvertes.fr/docs)

## Some libraries

- https://www.bibliometrix.org/home/
- https://aurelien-goutsmedt.com/post/extracting-biblio-data-1/
- https://github.com/sbegueria/bibliometRics
- searchlitr package r
- https://rdrr.io/github/nfrerebeau/odyssey/f/README.md
- Fulltext package

## Focus on crossrefs

rcrossref : le plus important ? mais fermé à la liste bibliographique cite-by.
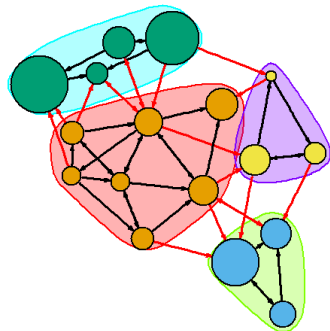
Section 3

Network analysis

Figure 5: A network and communities

- *igraph* the perfect tool with r and (an excellent introduction)[https://kateto.net/netscix2016.html]
- data : $x < -w_i- > y$
  - ▶ co-occurences and others distances.
- Analytical tools :
  - ▶ Layout ( MDS, KR, …)
  - ▶ Centrality measurements ( HITS, …)
  - ▶ Cliques and communities detection

# PMP Authorship Analysis

- PMP case : 40 years of publications - around 1020 papers
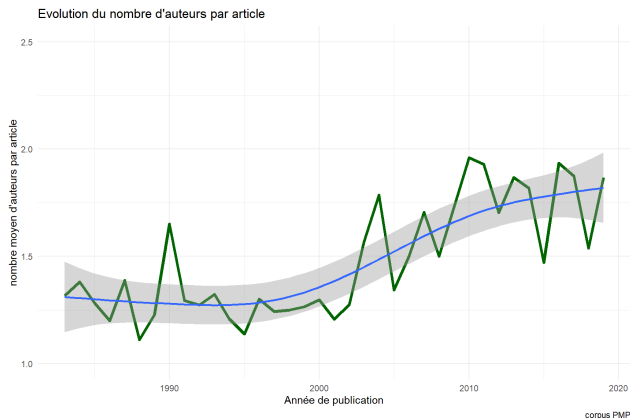- Tracking a regime change

Evolution du nombre d'auteurs par article



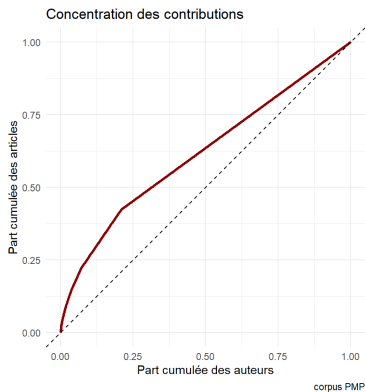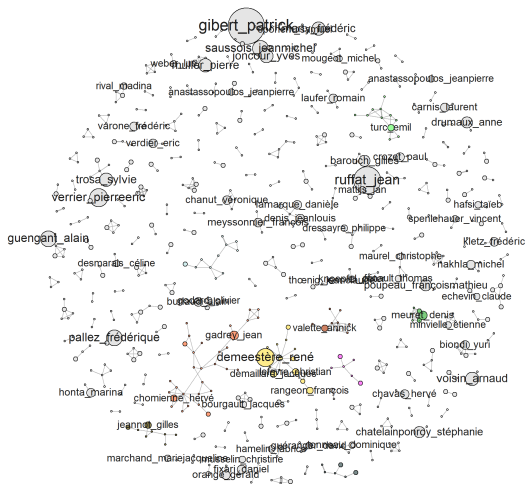Figure 6: Figure 10 : a change in authorship

Figure 7: Author concentration

**Les constellations des auteurs de PMP**



Sélection des auteurs ayant publié au moins une fois à plusieurs :
taille = nombre d'articles publiés
couleur : vert 2 et moins, coral : 3 et plus

Figure 8: network of authors

Projection Tsne des auteurs

# L'exemple d'un réseau de citation

Un petit exemple par olivier

Section 4

Topic Models and embeddings

# Topic Models and embeddings

- Topic model represent the first modern wave of text statistical modeling approach with LDA models Blei, Ng, and Jordan (2003).
- Embeddings a second wave with Mikolov et al. (2013)
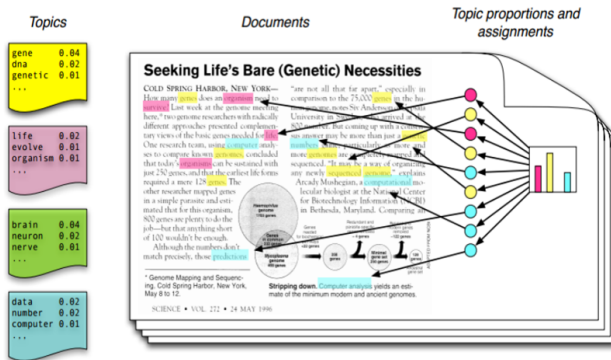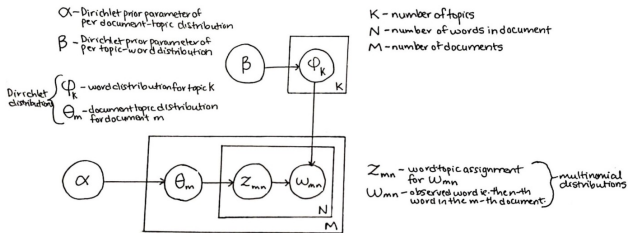- Tranformers is all you need, now.

# The LDA model



Figure 9: The LDA concept

α – Dirichlet prior parameter of per document-topic distribution

β – Dirichlet prior parameter of per topic-word distribution

Dirichlet distribution
- $φ_k$ – word distribution for topic k
- $θ_m$ – document-topic distribution for document m

K – number of topics
N – number of words in document
M – number of documents

$Z_{mn}$ – word-topic assignment for $W_{mn}$
$W_{mn}$ – observed word ie. the n-th word in the m-th document
} multinomial distributions

Now a large family with Structural Topic Models or Seed LDA models.

# An application of a STM model



**Top Topics**

- Topic 4: france, publics, services
- Topic 3: publique, l'action, fonction
- Topic 1: public, management, service
- Topic 6: politiques, publiques, analyse
- Topic 9: gestion, ressources, pays
- Topic 2: contrôle, développement, organisations
- Topic 5: d'un, l'exemple, face
- Topic 8: entre, vers, l'administration
- Topic 10: politique, d'une, recherche
- Topic 7: système, collectivités, performance
- Topic 11: locales, européenne, responsabilité
- Topic 12: cas, nouvelle, comment

Expected Topic Proportions
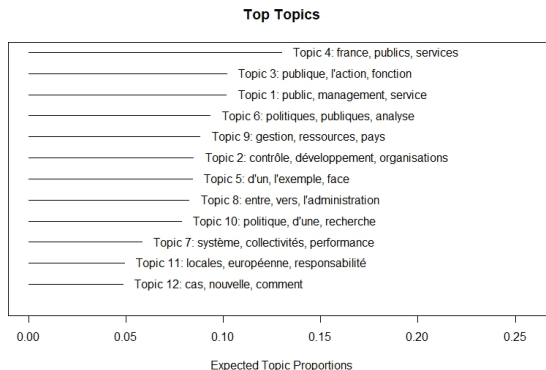
0.00    0.05    0.10    0.15    0.20    0.25

Figure 10: A STM topic model : keywords and proportion of the content /n (each document has a p probability to belong to the topic k)
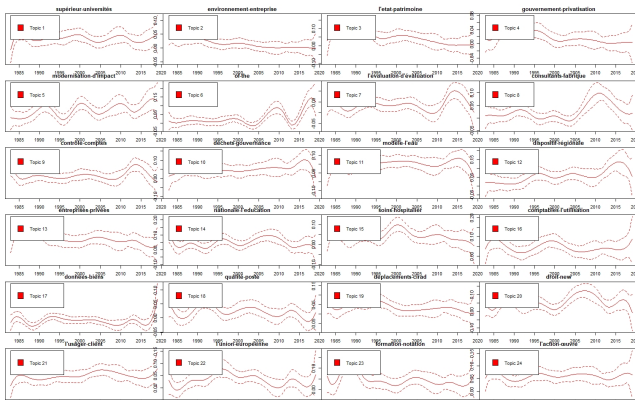
Figure 11: A STM topic model : time prevalence for each topic identified
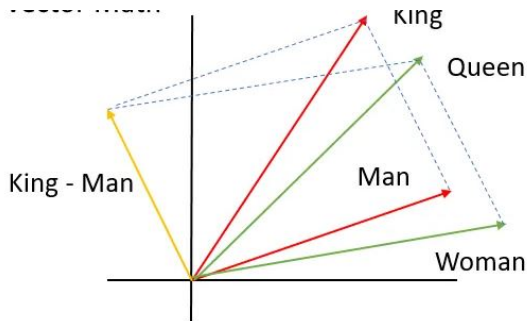
# An embeddings approach
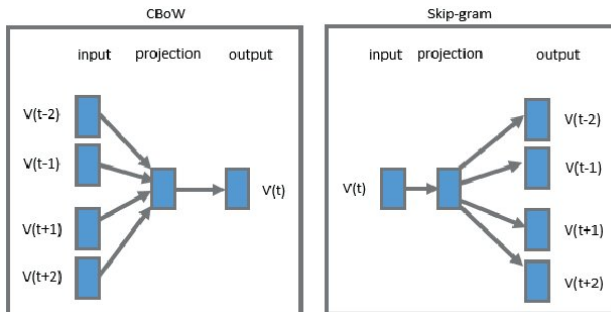


Figure 12: Embeddings intuition

Figure 13: The word2vec model

### Annotation stage

```
2023-02-16 08:33:40 Annotating text fragment 1/110
2023-02-16 08:33:46 Annotating text fragment 21/110
2023-02-16 08:33:54 Annotating text fragment 41/110
2023-02-16 08:34:02 Annotating text fragment 61/110
2023-02-16 08:34:09 Annotating text fragment 81/110
2023-02-16 08:34:18 Annotating text fragment 101/110
Time difference of 42.73675 secs
```
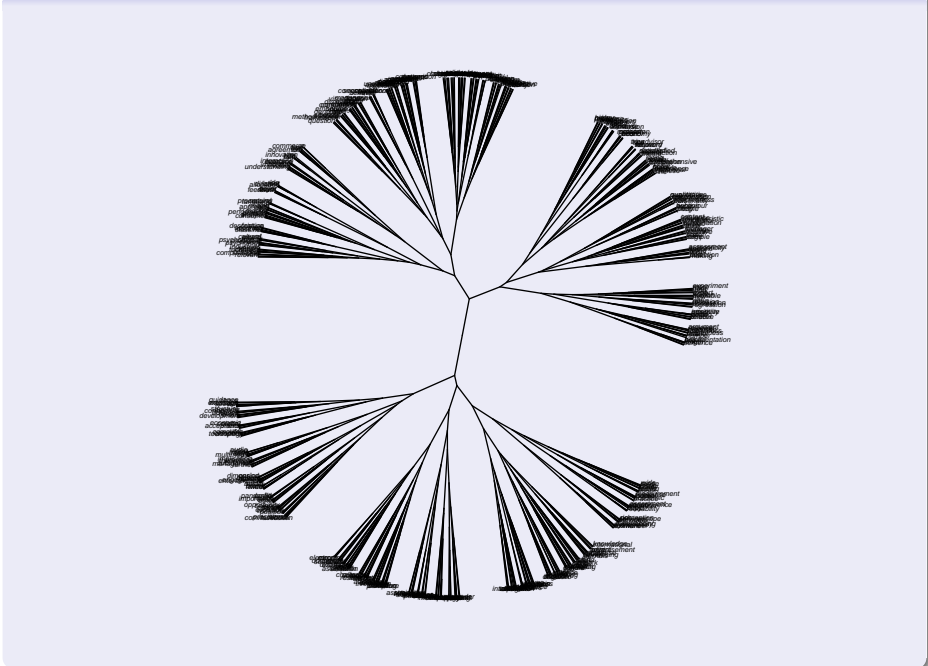
## vectorisation Stage

Un test

```
$lda
    term1         term2 similarity rank
1     lda      dirichlet 0.8190907    1
2     lda     allocation 0.8050508    2
3     lda         latent 0.6876944    3
4     lda        mindset 0.6287284    4
5     lda  psychological 0.6183707    5
6     lda         course 0.6016998    6
7     lda          award 0.5977740    7
8     lda    description 0.5879406    8
9     lda      awareness 0.5843993    9
10    lda        chinese 0.5599142   10
11    lda          model 0.5524023   11
12    lda          topic 0.5437085   12
13    lda        outcome 0.5366894   13
14    lda          other 0.5325939   14
15    lda    methodology 0.5320520   15
16    lda        balance 0.5313188   16
17    lda     industries 0.5309694   17
18    lda       relevant 0.5299399   18
19    lda           loan 0.5275514   19
20    lda         seller 0.5235679   20
```
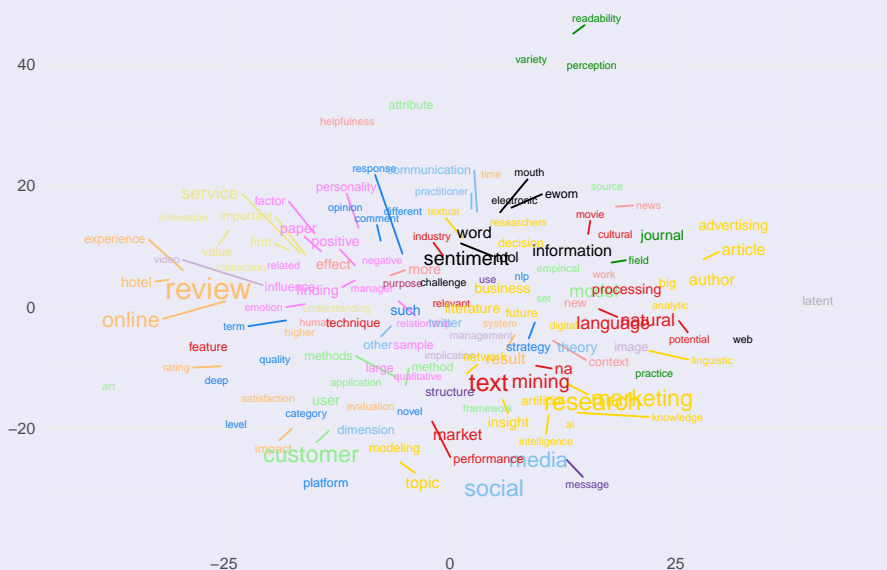
1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16

Projection 2D du vocabulaire vectorisé des 110 articles 'NLP & Marketing'
Pipe : annot. syntax –> word2vec: 200 vecteurs–> hclus –>rtsne
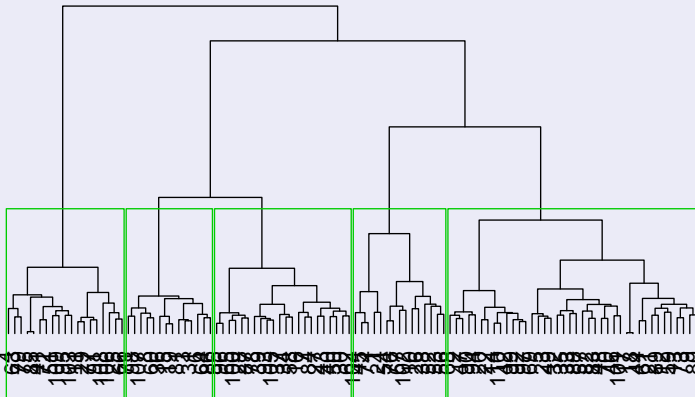
### document embeddings

Text vectors are created simply by calculating the resulting vector of words in the text :
text $= le\ NLP\ quantifie\ les\ mots$
$V_text = V_n lp + V_q uantifier + V_m ots$
We could also compute abstract concepts. for exemple, finding mention of advanced NLP methods : - concept="LDA Word2vec Bert ML" - compute the similarity between each text and the concept

grouping documents

**Cluster Dendrogram**

document projection space

Projection 2D des 106 articles 'NLP & Marketing' vectorisés
Pipe : annot. syntax –> word2vec: 200 vecteurs–> doc2vec–> hclus –>rtsne

# We might as well put everything in one space

Because embedding space is common.

Projection 2D du vocabulaire vectorisé des 110 articles 'NLP & Marketing'
Pipe : annot. syntax –> word2vec: 200 vecteurs–> hclus –>rtsne

Section 5

Conclusion

# Conclusion

- Another organisation of the literature review that requires a collective effort.
  - building common code
  - discussion arena
- Toward a systematic understanding of bibliometrics metadata
  - intrinsic : content and localization (Journal, date, institution)
    - first order information : pure meta data
    - second order information : outcomes, methods,
  - extrinsic : how papers related to the paper world
    - citation index,
    - co-references
- Be prepared for the disruption of deep NLP methods :
  - summarizing
  - concept extraction
  - outcome extraction
  - and more...

Section 6

References

# References

Blei, David M., Andrew Y. Ng, and Michael I. Jordan. 2003. "Latent Dirichlet Allocation."
  *J. Mach. Learn. Res.* 3 (March): 993–1022.
  http://dl.acm.org/citation.cfm?id=944919.944937.
Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. "Efficient Estimation
  of Word Representations in Vector Space." arXiv. http://arxiv.org/abs/1301.3781.
Moher, D., A. Liberati, J. Tetzlaff, D. G Altman, and for the PRISMA Group. 2009.
  "Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA
  Statement." *BMJ* 339 (July): b2535–35. https://doi.org/10.1136/bmj.b2535.
Stratton, Samuel J. 2019. "Literature Reviews: Methods and Applications." *Prehospital
  and Disaster Medicine* 34 (4): 347–49. https://doi.org/10.1017/S1049023X19004588.