

# Grado en Inteligencia Artificial

## Textual Analysis – CA Approach

**Prof. Dante Conti**  
**Prof. Sergi Ramírez**

- ① Introduction
- ② Methods and approaches
- ③ Correspondence Analysis applied to textual data
- ④ Example 1: Definition of health
- ⑤ Example 2: Comparison of factor analysis methods
- ⑥ Summary and conclusions

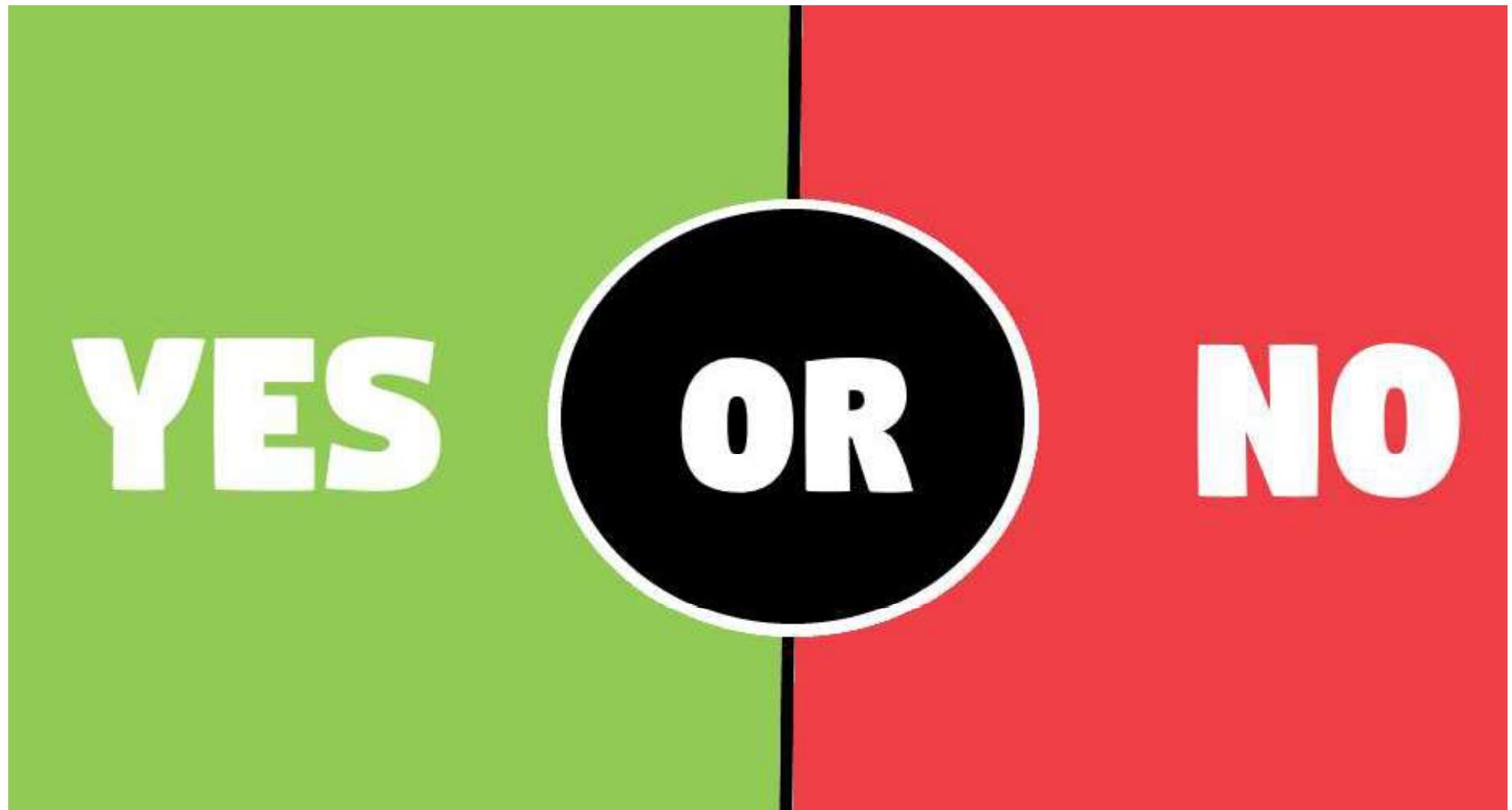


Imagine that tomorrow you  
will go to the Liceu to see “La  
Traviata” with your best friend

...



And at the end, you ask  
your friend:  
- “Did you like it?”





And at the end, you ask  
your friend “What is your  
opinion about it?”

# Introduction





deceit kill deep flaw blame haunt tension weird bad  
crap betray loss fake complex  
strike junk spite shatter  
stun shit flee disappoint freak  
die twist mean reject plot killer hung  
bore fall overwhelm low suck fuck  
steal lack sadimpromptu bleed tedious  
difficult disown stuck harass  
shock pleas critic slap pain  
slow miser daze small scream  
stole poor beat less  
miss hell

What do you expect when you go on holidays?



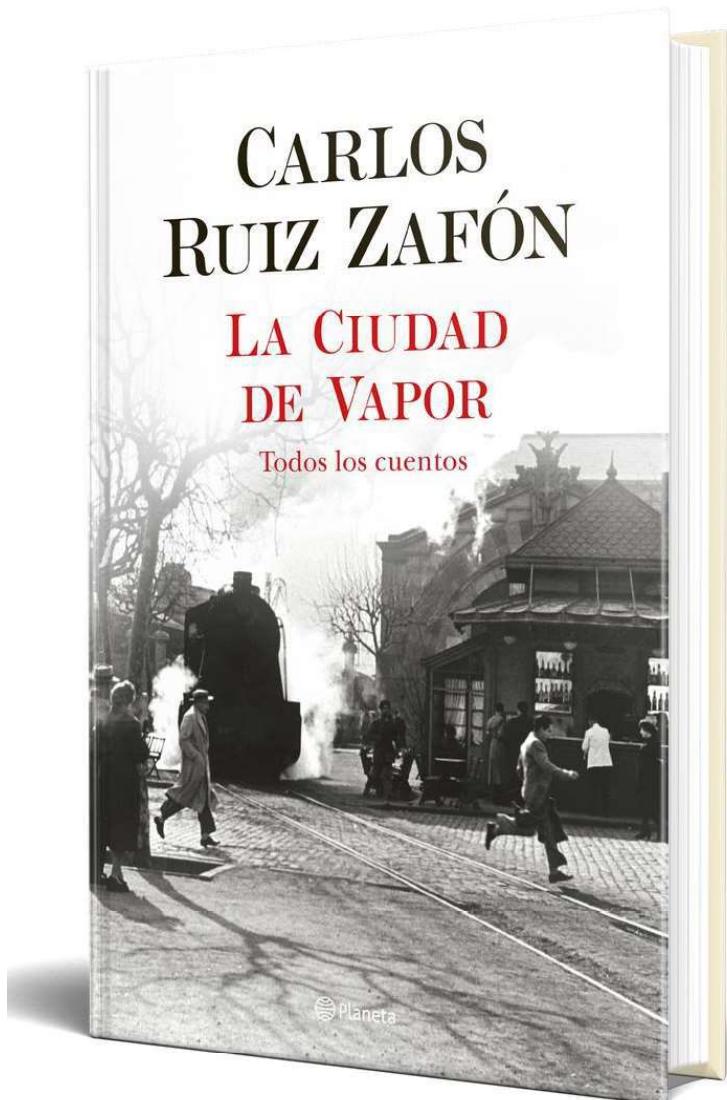
UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH



# What do you expect when you go on holidays?

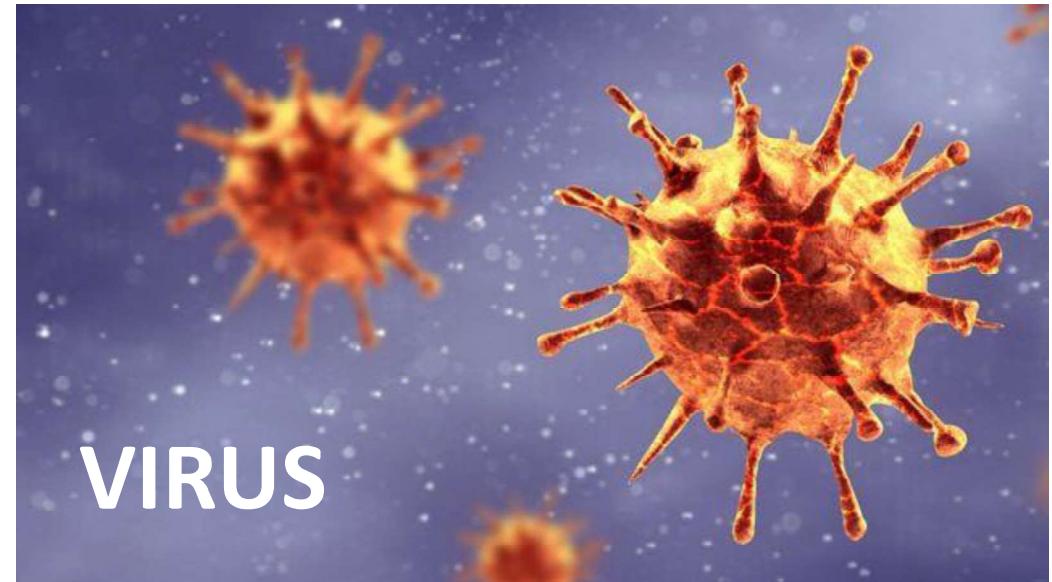


# WORDS



«Puedo conjurar rostros de chiquillos del barrio de la Ribera con los que a veces jugaba o peleaba en la calle, pero ninguno que quisiera rescatar del país de la indiferencia. Ninguno excepto el de Blanca.»

What comes to your mind when you think about “COVID-19”?



## I Have a Dream (delivered 28 August 1963, at the Lincoln Memorial, Washington D.C.)

I am happy to join with you today in what will go down in history as the greatest demonstration for freedom in the history of our nation.

Five score years ago, a great American, in whose symbolic shadow we stand today, signed the Emancipation Proclamation. This momentous decree came as a great beacon light of hope to millions of Negro slaves who had been seared in the flames of withering injustice. It came as a joyous daybreak to end the long night of their captivity.

But one hundred years later, the Negro still is not free.....





**Donald J. Trump** 

@realDonaldTrump



 Follow

Healthy young child goes to doctor, gets pumped with massive shot of many vaccines, doesn't feel good and changes - AUTISM. Many such cases!

RETWEETS

8,768

LIKES

6,283



8:35 AM - 28 Mar 2014

Deliver to  Spain All ▾ electric scooter

Deals Customer Service Gift Cards Sell Registry

of 253 results for "electric scooter"

Sort by: Relevance

Department: Scooters

Sub Departments: Kick Scooters, Balance Scooters, Other Parts, Scooter Batteries & Battery Chargers, E-Bikes, All 10 Departments

Customer Review:    

Price: Under \$25, \$25 to \$50, \$50 to \$100, \$100 to \$200, \$200 & Above

Shop by price: Price and other details may vary based on size and color

  
Sponsored   
Gotrax GXL V2 Commuting Electric Scooter - 8.5" Air Filled Tires - 15.5MPH & 9-12 Mile Range - Version 2  
 ~ 2,773  
**\$248<sup>00</sup>** ~~\$299.99~~  
Ships to Spain

  
Sponsored   
Gotrax XR Ultra Electric Scooter, LG Battery 36V/7.0AH Up to 17 Miles Long-Range, Powerful 300W Motor & 15.5 MPH, UL Certified Adult E-Scooter  
 ~ 496  
**\$329<sup>99</sup>** ~~\$399.99~~  
Ships to Spain

  
Sponsored   
Mongoose React Electric Kids Scooter, Boys & Girls Ages 8+, Max Rider Weight Up To 175lbs, Varying Max Speed, Aluminum Handlebar  
 ~ 31  
**\$119<sup>99</sup>**  
Ships to Spain

# Customer reviews

### Customer reviews

★★★★★ 4 out of 5

2,773 global ratings

Rating	Percentage
5 star	61%
4 star	14%
3 star	6%
2 star	5%
1 star	15%

▼ How are ratings calculated?

### By feature

Feature	Rating
Maneuverability	★★★★★ 4.2
Warmth	★★★★★ 3.8
For commuting	★★★★★ 3.7

▼ See more

### Review this product

Share your thoughts with other customers

Write a customer review

  
Crisp Inception Pro Scooter  
\$119.99 prime  
[Shop now](#)

Sponsored

### Customer images



See all customer images

### Read reviews that mention

battery life    cruise control    top speed    stopped working    full charge  
front wheel    fully charged    year old    great scooter    day warranty  
front tire    make sure    highly recommend    put together    air filled

Most recent ▾

### From the United States

 Robinson C.  
★★★★★ Great  
Reviewed in the United States on December 1, 2020  
Color: Black | Verified Purchase  
Good to move around the neighborhood  
[Helpful](#) | [Comment](#) | [Report abuse](#)

 Alvaro  
★★★★★ Fast scooter and fast shipping!  
Reviewed in the United States on December 1, 2020  
Color: Black | Verified Purchase  
This is your go-to electric scooter. For the price, the power and the range this scooter has beats any scooter even top brands, i had to buy a second one so my girl could come. Im reviewing a 5 month old scooter i have used it everyday on late night rides and whenever i go out to ride with family, etc. Both of them work great and havent had any problems and i only have positive things to say about it i havent found a flaw on it. Maybe the battery indicator its not the most accurate but other than that is great chargest fast and it is fast!!  
[Helpful](#) | [Comment](#) | [Report abuse](#)

 Jennifer A.  
★★★★★ Box is cut open  
Reviewed in the United States on November 29, 2020  
Color: Black | Verified Purchase  
Box is cut open and a bit mangled and we hope the scooter works good.  
One person found this helpful

- To bring up the subject
- To determine the structure of the text:  
localizing the ruptures
- To visualize the proximity between documents  
and words
- To index the documents
- To construct a database with documents for a  
subsequent automatic interrogation

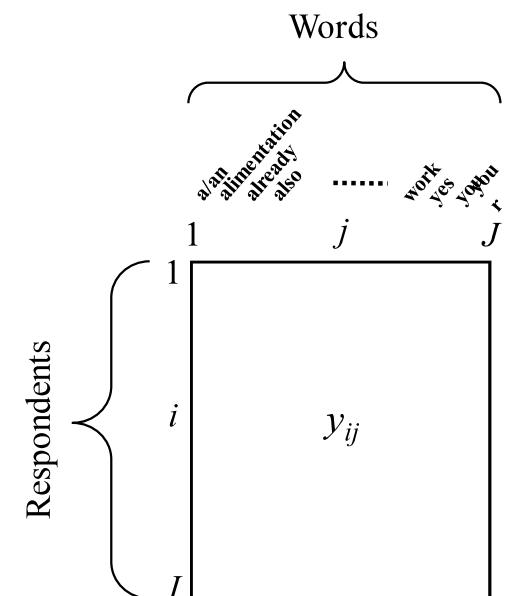
- Classical methods: frequency of words
- Multidimensional methods: Factor analysis
- Recent methods: Machine learning
- Two important points: codification and distance between the elements

- Surveys often include open-ended questions as a complement to classical closed questions to address complex and little known topics.
- Example: survey intended to better know the definitions of health that the non-experts give
  - 392 respondents who answered through free-text comments
  - The respondents' characteristics: age in groups (under 21, 21–35, 36–50 and over 50), gender (man and woman) and health condition (poor, fair, good and very good health)

To obtain the final word frequency table to analyze, we :

- Replace special characters from the text (for example, replacing “/”, “@” and “|” with space).
- Remove unnecessary white space
- Convert the text to lower case
- Convert the verbs to infinitive form
- Remove common stopwords like “the”, “we”, etc.

- The  $(I \times J)$  matrix  $Y$ , with generic term  $y_{ij}$ , contains the frequency of the  $J$  words in the  $I$  respondents' answers.
- Stopwords and lemmatization
- The minimum threshold for words: 10
- 115 different words and 7751 occurrences ( $N$ )



- Correspondence Analysis (CA) is the reference factorial method
  - Optimal visualization of the similarities among the free answers
  - Optimal visualization of the similarities among the words
  - The associations between words and free answers.

- Correspondence analysis (CA) was proposed by *Benzecri* (1973, 1981) to deal with textual data.



- The term "correspondence analysis" stems from the fact that lexical table links two corresponding sets:
  - one represented by the rows, and
  - other represented by the columns, playing symmetric roles.

- The  $(I \times J)$  frequency table  $\mathbf{Y}$  is transformed into a proportion table

$$\mathbf{P} = \mathbf{Y}/N$$

with generic term  $p_{ij} = y_{ij}/N$ .

- The  $(I \times I)$  diagonal matrix  $\mathbf{D}$  stores the row margins  $p_{i\bullet} = \sum_{j \in J} p_{ij}$
- The  $(J \times J)$  diagonal matrix  $\mathbf{M}$  stores the column margins  $p_{\bullet j} = \sum_{i \in I} p_{ij}$

- Classical CA results can be obtained through a factor analysis applied to the ( $I \times J$ ) matrix

$$X = D^{-1}PM^{-1} = [q_{ij}] = [p_{ij}/p_{i\bullet}p_{\bullet j}]$$

with metrics/weights **M** and **D** (*Escofier and Pages, 1988*).

- The inertia axes in the row space satisfy the equation  $\mathbf{X}'\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{U} = \mathbf{U}\Lambda$
- The vector of coordinates of the rows  $\mathbf{F} = \mathbf{X}\mathbf{M}\mathbf{U}$
- The inertia axes in the column space satisfy the equation  $\mathbf{X}\mathbf{M}\mathbf{X}'\mathbf{D}\mathbf{V} = \mathbf{V}\Lambda$
- The vector of coordinates of the columns  $\mathbf{G} = \mathbf{X}'\mathbf{D}\mathbf{V}$
- The simultaneous representation of rows and columns relies on the transition relationships

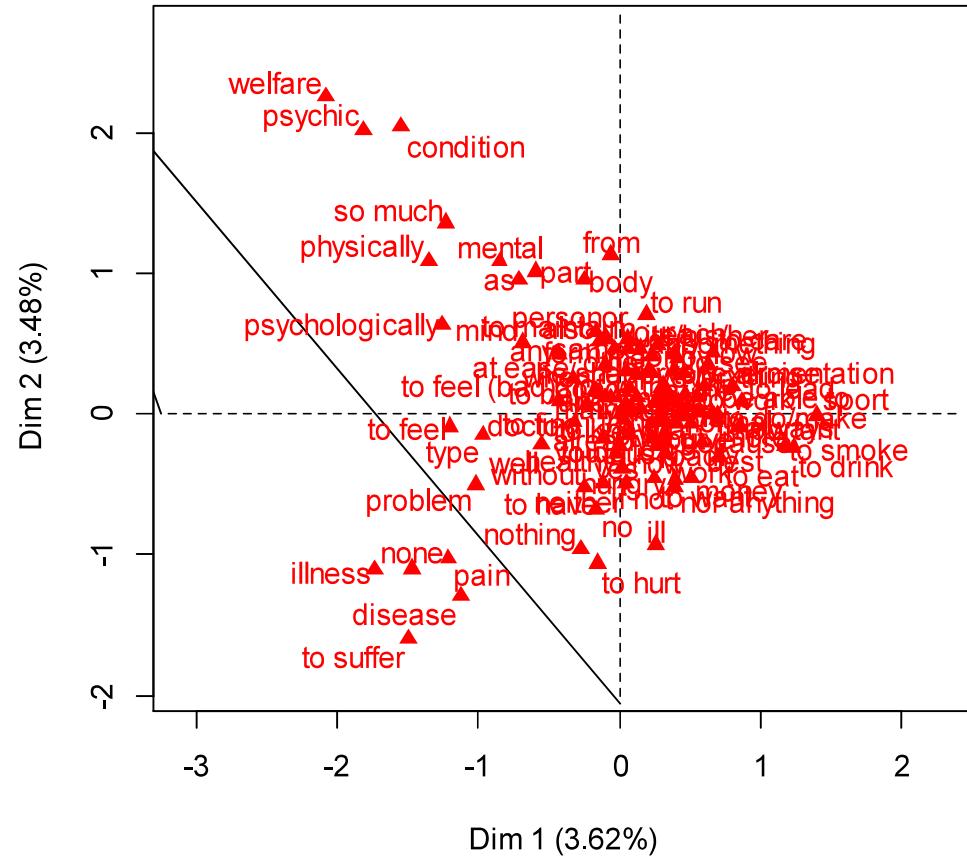
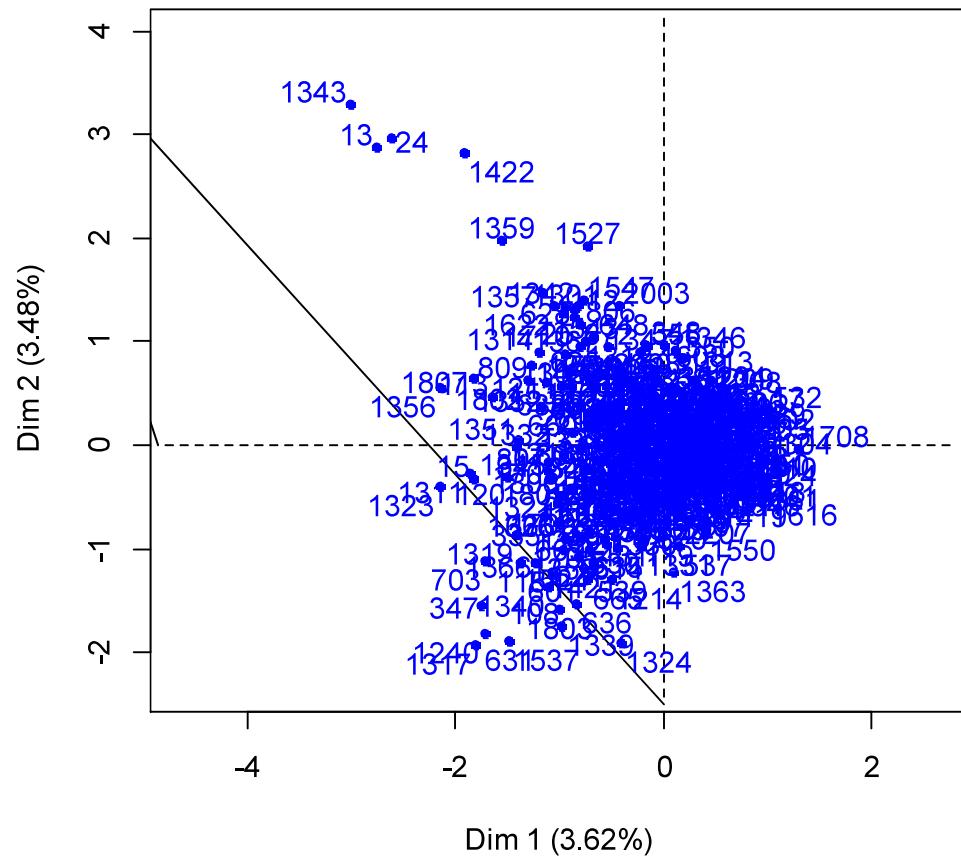
$$\mathbf{F} = \mathbf{V}\Lambda^{1/2}$$

$$\mathbf{G} = \mathbf{U}\Lambda^{1/2}$$

- For dimension  $s$  of this superimposed representation, up to the multiplicative factor  $1/\sqrt{\lambda_s}$ , a row  $i$  (resp. a column  $j$ ) is at the barycentre of the columns (resp. of the rows), each column  $j$  (resp. each row  $i$ ) having a weight  $p_{ij}/p_{i\bullet}$  (resp.  $p_{ij}/p_{\bullet j}$ ). This property is used to interpret the position of one row in relation to all of the columns and the position of one column in relation to all of the rows.
- The similarity between rows and columns is expressed in a totally symmetrical way (characteristic of CA which differentiates it from other principal components methods). Two rows are all the closer as they are frequently associated with the same columns and two columns are all the closer that they are frequently associated with the same rows.

	eigenvalue	percentage of variance	cumulative percentage of variance
dim 1	0.2333100	3.619220	3.619220
dim 2	0.2243156	3.479693	7.098914
dim 3	0.1858920	2.883649	9.982563
dim 4	0.1753531	2.720163	12.702726
dim 5	0.1502916	2.331397	15.034123
dim 6	0.1433485	2.223692	17.257814
dim 7	0.1418055	2.199757	19.457572
dim 8	0.1348692	2.092157	21.549729
dim 9	0.1325759	2.056582	23.606311
dim 10	0.1293922	2.007196	25.613506
		.....	
dim 114	0.0070151	0.108822	100.000000

# Results CA: first factorial plane



	<b>Dim 1</b>	<b>Dim 2</b>
physically	10.4460384	6.94756938
none	10.9328585	6.37469296
welfare	5.0371786	6.15360219
no	0.5879697	10.53125665
disease	4.5966688	6.36282802
psychic	4.1886973	5.36576630
condition	3.3205504	5.97001380
to feel	8.6498820	0.04843705
to have	0.3388698	6.13595298
nothing	0.3892418	4.96543895

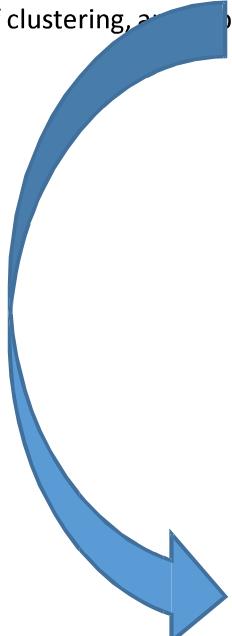
# From text to a word frequency table

**PCA.** This week, we have for you three videos which together present the main details of principal component analysis. Principal component analysis is a set of tools which allow us to study and visualize large data sets. We will present the method from a theoretical as well as a practical point of view....

**CA.** This week, we have for you 5 course videos on correspondence analysis. In the videos, we will see the following: first, we start by describing the data, giving a little notation, and considering questions to ask when running correspondence analysis. We'll see that the main point of correspondence analysis is studying the links between pairs of qualitative variables....

**MCA.** This week, we have four videos for you on multiple correspondence analysis, MCA for short. We'll have a look at the main features of the method, using a specific example to guide us along the way. The videos look at the following things. First, we describe the types of data MCA can be used for. With this data in mind, we will look at what our goals are, and what issues we may have....

**Clustering.** This week, we're going to look at clustering methods, including hierarchical clustering, and a partitioning method called k-means. The course videos for this week get into the following things: After a brief introduction on data of interest for clustering, and the goals of clustering, we are going to have a look at some general principles of clustering, and in particular, hierarchical ...



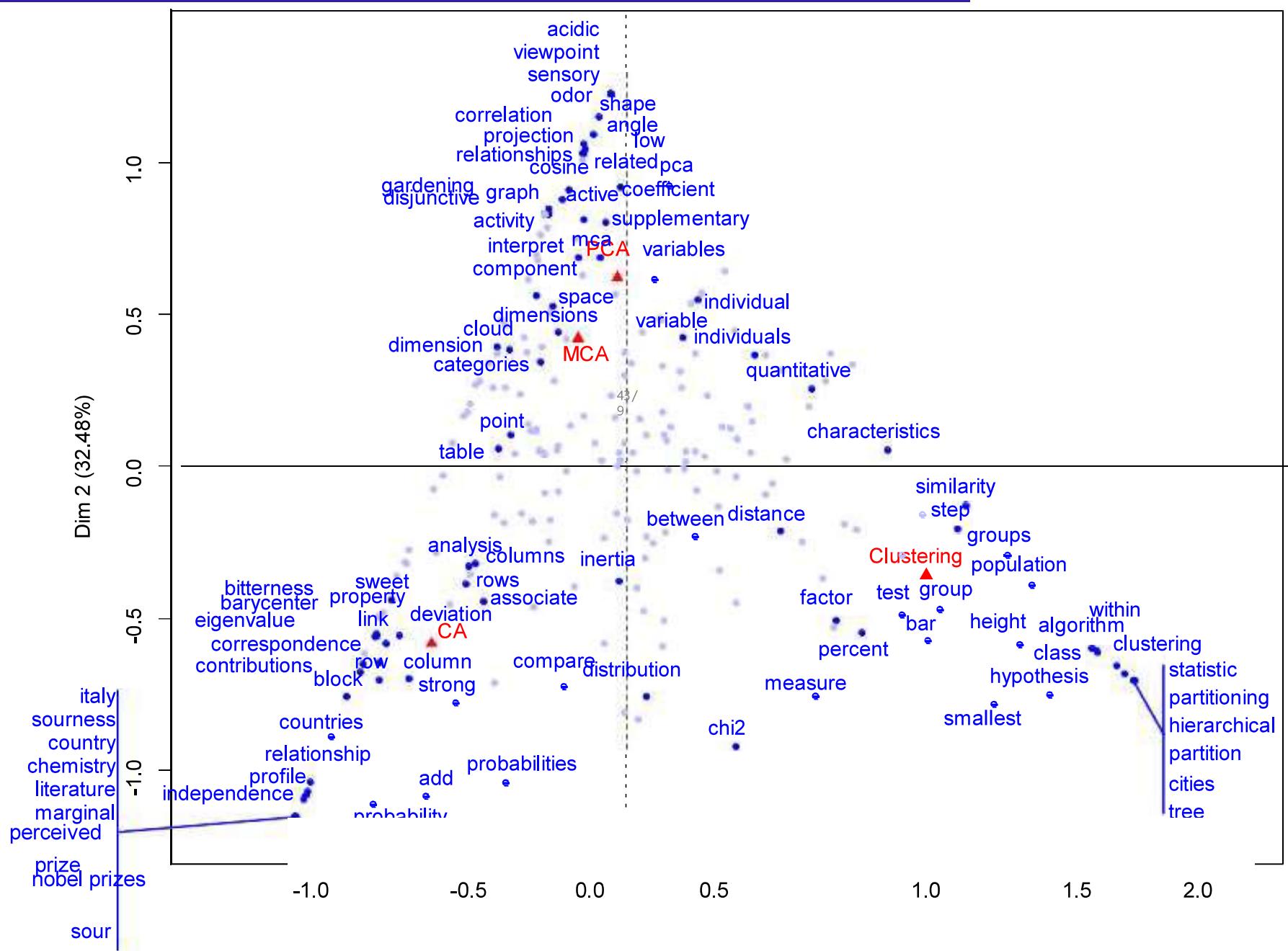
1852 rows

	PCA	CA	MCA	Clustering
able	2	1	1	2
above	0	0	0	1
absolute	1	0	0	5
absolutely	0	1	0	0
acceptable	0	0	0	1
access	1	0	0	0
accident	0	0	1	0
accord	0	0	1	0
according	3	0	2	0
account	0	0	2	1
...	...	...	...	...

- remove connecting words like : for example, then, therefore, and, etc.
- group words with the same root or the same conjugations together (e.g., reduced, reduction, reduces)
- group singular and plurals together
- remove words used nine times or less

⇒ 246 words,  $n = 8821$  total occurrences

# Simultaneous representation of methods and words



Terms exclusive to certain methods are superposed

1st axis (inertia = 0.35) :

clear division between the factor methods and clustering

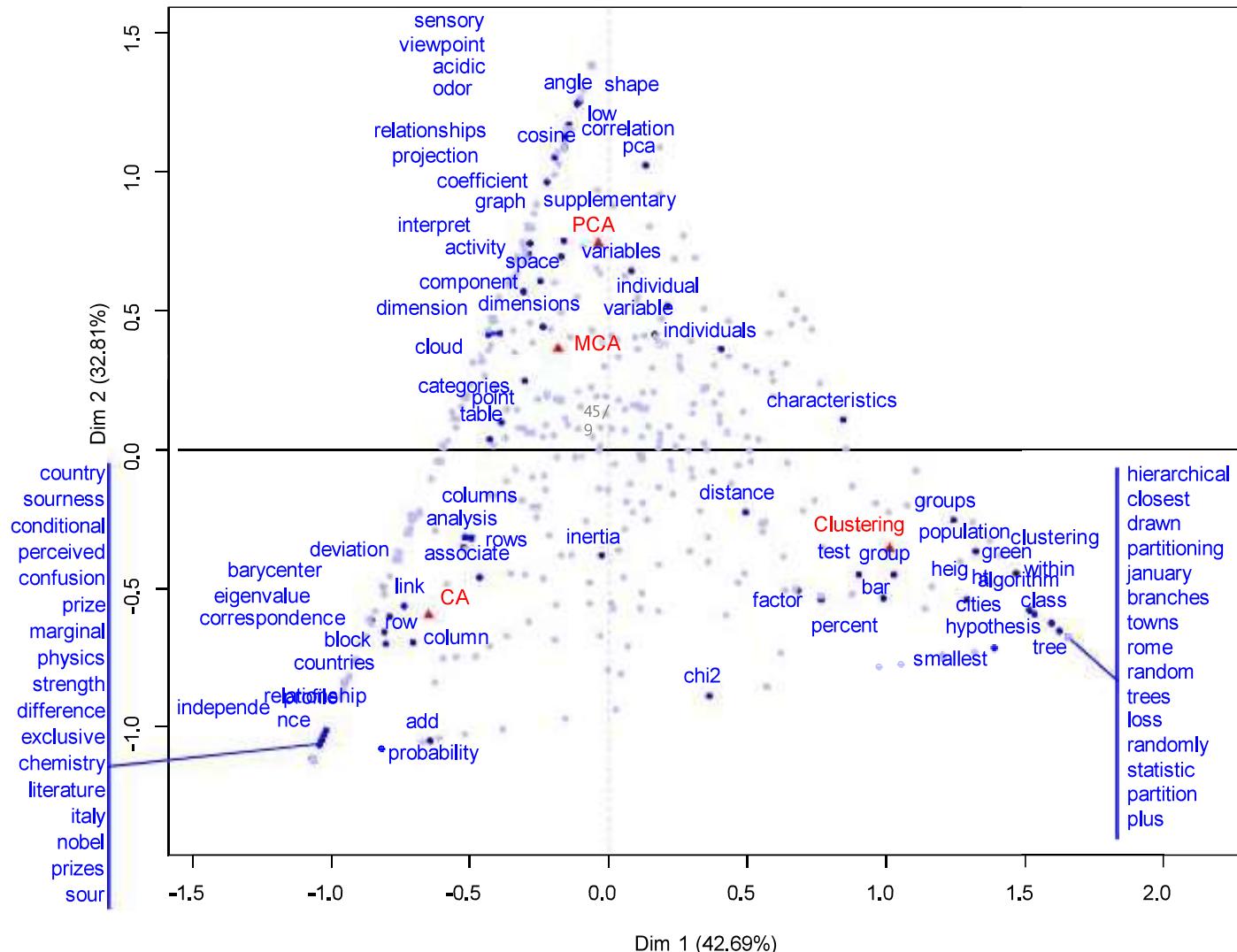
specific words used in factor analysis are to the left

specific words used in clustering are to the right

2nd axis (inertia = 0.26) :

separates the 3 factor analysis methods

## Words included : used at least 5 times



==> Stable representation

For more information about CA applied on textual data, have a look at this video:

<https://www.youtube.com/watch?v=62N0JMf5hOs>

**Husson F., Lê S. & Pagès J.** (2017) *Exploratory Multivariate Analysis by Example Using R* 2nd edition, 230 p., CRC/Press.

**Benzecri, J. P.** (1973). Analyse des Donnees. Dunod.

**Benzecri, J. P.** (1981). Pratique de l'analyse des donnees: Linguistique & lexicologie, volume 3. Dunod.

**Escofier, B. and Pages, J.** (1988). Analyses factorielles simples et multiples. Dunod.

**Lebart, L., Salem, A., and Berry, L.** (1998). Exploring Textual Data. Kluwer Academic Publishers.