

# Clustering de textes courts pour relier des députés

Présentation de stage technique, MAIN 4

Achille Baucher

Encadré par Marie-Jeanne Lesot et Adrien Revault d'Allones

Octobre 2020



# Présentation

## Données

- 20 000 lettres
- ~ 160 mots

**Député** : Bertrand Sorre

**Ministre** : Ministre de l'intérieur

**Date** : 12 Juin 2018

**Rubrique** : Sécurité routière

**Titre** : Précision sur les 80 km/h

**Lettre** :

M. Bertrand Sorre attire l'attention de M. le ministre d'État, ministre de l'intérieur, **sur la décision prise de réduire la vitesse maximale autorisée à 80 km/h**, à compter du 1er juillet 2018, sur les routes à double sens sans séparateur central (limitée actuellement à 90 km/h). Cette disposition permettra de sauver entre 300 et 400 vies par an selon le comité des experts du conseil national de la sécurité routière dans son rapport du 29 novembre 2013. Fréquemment questionné à ce sujet dans la circonscription de La Manche dont il est l'élu, *il aimerait savoir si cette décision implique également des modifications sur les vitesses actuellement autorisées pour les professionnels de la route (transport routier, autobus) ou pour les apprentis conducteurs d'un véhicule léger, titulaires d'un permis depuis de moins de 2 ans.*

# Présentation

## Objectif

- Clustering des lettres
- Même sujet précis
- Lien entre députés



# Présentation

## Normalisation

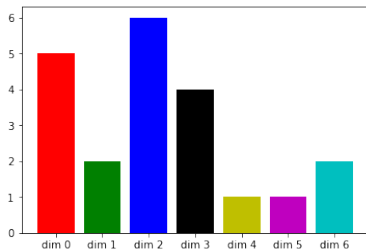
*M. Bertrand Sorre attire l'attention de M. le ministre d'État, ministre de l'intérieur, sur la décision prise de réduire la vitesse maximale autorisée à 80 km/h.*

- Ponctuation
- Stopwords
- Minuscule
- Infinitif singulier

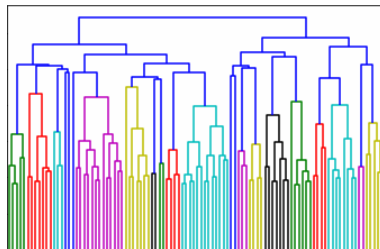
→ *bertrand sorre attirer attention ministre état ministre intérieur décision prendre réduire vitesse maximale autorisation 80 km/h.*

# Clustering de documents courts

## État de l'art



• Représentation vectorielle



• Clustering classique

# Représenter un document sous forme de vecteur

## Principe et TF

- Un mot par dimension
- Term-Frequency

Documents \ Mots	Covid	chloroquine	...	santé
Document 1 (TF)	4	2	...	3
Document 2 (TF)	2	0	...	3

# Représenter un document sous forme de vecteur

## TF-IDF

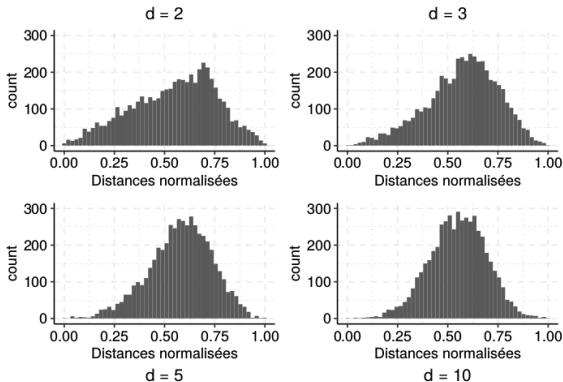
- Term-Frequency Inverse Document Frequency
- + Mots singuliers
- - Mots communs

Documents \ Mots	Covid	chloroquine	...	santé
Document 1 (TF)	4	2	...	3
Document 2 (TF)	2	0	...	3
Document 1 (TF-IDF)	3.	3.5	...	1.2
Document 2 (TF-IDF)	1.5	0.	...	1.2

- $n \approx 9000$  dimensions !

# Malédiction de la dimensionnalité

Exemple : Concentration des distances



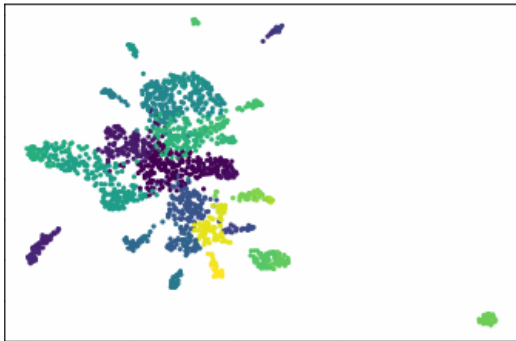
- Clustering difficile



# Réduire la dimension

## Enjeux

- Réduction drastique
- Préserver l'information
- Repérer la structure
- Diverses méthodes



# Réduire la dimension

## Principe

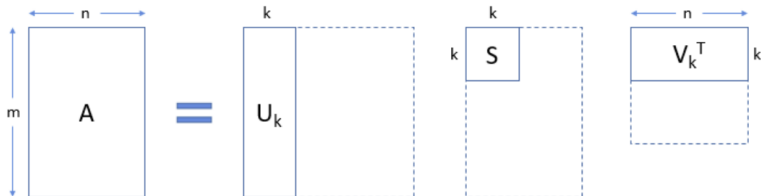
- Partir des TF ou TF-IDF
- Choisir P
- Nouvelles dimensions ( $\sim$  *thèmes*)

Documents	dim 1 ( <i>Ex : Santé</i> )	dim 2 ( <i>Urgence</i> )	...	dim P ( <i>BTP</i> )
...blabla Covid...	3.5	4.7	...	0.0
...blabla EDF ...	0.0	0.6	...	3.7

# Réduire la dimension

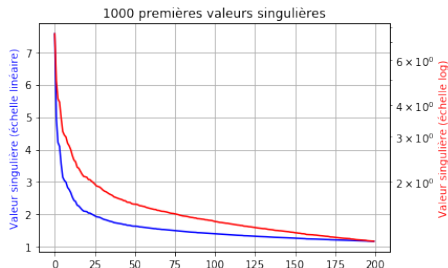
## Algorithmes utilisés

- Factorisation du TFIDF (NMF, LSA)
- Distribution de probabilités des mots (LDA)

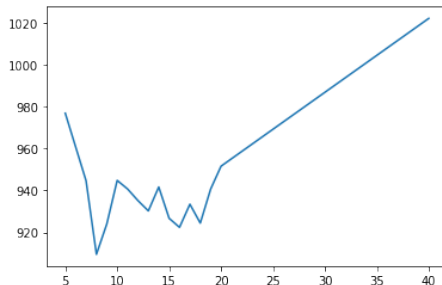


# Réduire la dimension

## Choix de la dimension



Critère de choix (LSA, NMF)



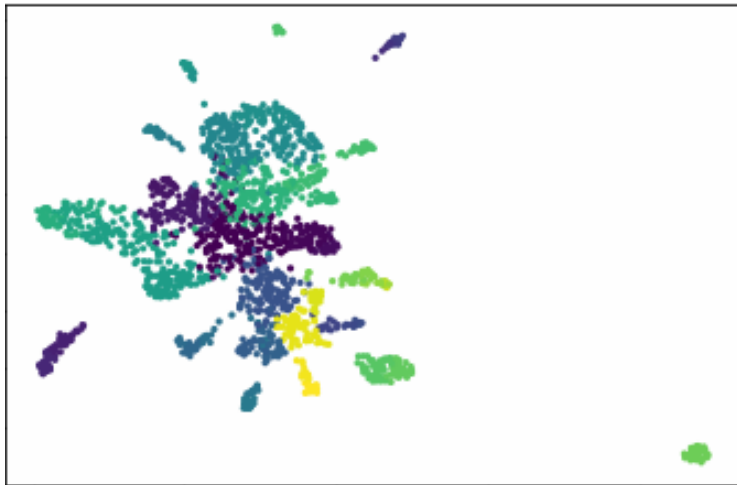
# Réduire la taille de dimension

## Conclusion

Vecteur	Topic 1	Topic 2	Topic 3
LSA	maladie, entreprise, santé, pouvoir, charge, mme, public, demande, prise, français	maladie, santé, lyme, prise, patient, charge, cancer, malade, atteindre, fibromyalgie	carte, gris, ant, us-ager, préfecture, service, titre, site, délai, sécuriser
LDA	tabac, cigarette, commune, logiciel, buraliste, tiers, prix, paquet, retrait, public	boulangerie, journée, administratif, paneterie, frontalier, tabac, effet, dispositif, gérant, hausse	art, métier, loi, roumanie, exercice, lme, fromager, oniam, examiner, lyme
NMF	000, pajemploi, pajot, pakistan, palais, palette, palier, palliatif, pallier, palmarès	charge, prise, patient, santé, atteindre, maladie, traitement, solidarité, autorité, soin	carte, gris, véhicule, problème, intérieur, demande, service, particulier, mois, agence

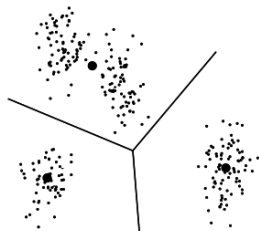
# Clustering

## Principe

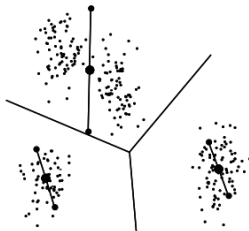


# Clustering

## K-means + X-means



K-means



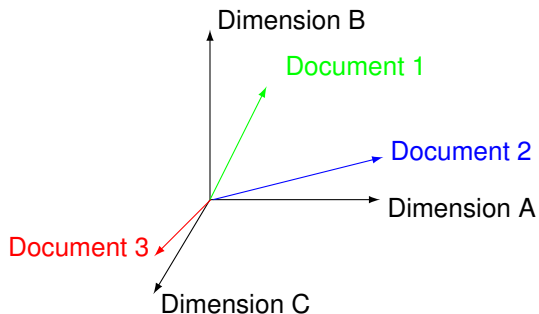
X-means step 1



X-means step 2

# Distances

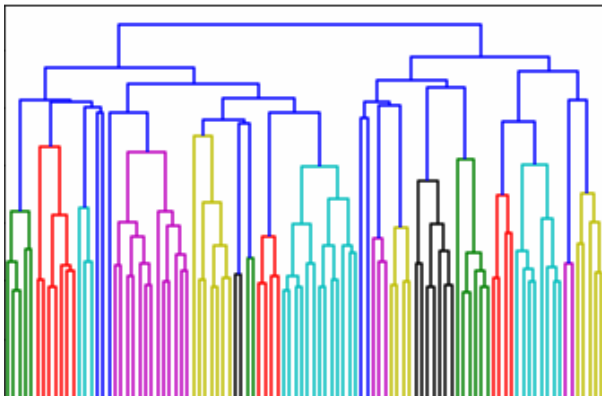
- Position ou direction ?





# Clustering

## Clustering Ascendant Hiérarchique (HAC)



# Clustering

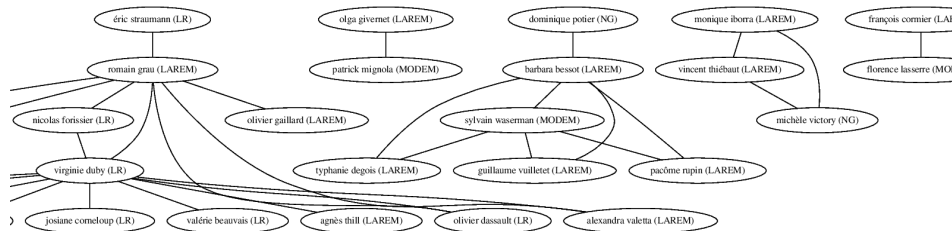
## Conclusion

- LSA et HAC
- NMF et X-means

Cluster A
déficit école national administration
déficit compte ena
ant - dématérialisation - dysfonctionnement péritoine
ena - déficit - gestion
...

# Lien entre députés

## Premier graphe



# Conclusion

