

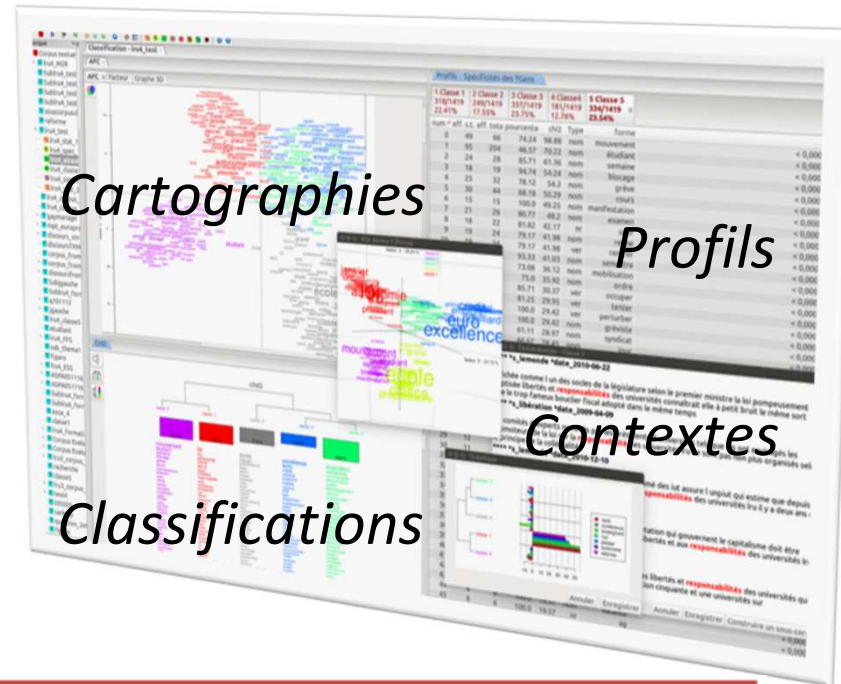
Journées d'étude @Media : cycle de l'information et Digital Methods.
Grenoble 30-31 mai 2018

L'analyse lexicométrique des discours en ligne :

méthodes et exemples d'application aux médias et espaces d'expression sur le web

Pascal MARCHAND, LERASS, Université Toulouse 3
Emmanuel MARTY, GRESEC, Université Grenoble Alpes

Analyser des corpus de textes



Cartographies

Profils

Contextes

Classifications

Iramuteq (<http://www.iramuteq.org>) est développé par Pierre Ratinaud au sein du *Lerass* et dans le cadre du *LabEx SMS*.



Plan de la présentation

- **Principes de base de la statistique lexicale**
Segmentation, partition, tableau lexical
- **Questions de recherche et résultats d'analyses: une typologie**
Mesurer et caractériser la diversité des expressions
Comparer plusieurs corpus
Caractériser un processus



Principes de base de la statistique lexicale

Segmentation et partition: deux opérations constitutives de la statistique lexicale

Quelques définitions:

Les questions que se donne la statistique lexicale sont les suivantes : « quels sont les textes les plus semblables en ce qui concerne le vocabulaire et la fréquence des formes utilisées ? Quelles sont les formes qui caractérisent chaque texte, par leur présence ou leur absence ? » (Lebart & Salem, 1994).

→ Tableau lexical (formes * textes)

La lexicométrie regroupe " toute une série de méthodes qui permettent d'opérer des ré-organisations formelles de la séquence textuelle et des analyses statistiques portant sur le vocabulaire à partir d'une segmentation " (Salem, 1986)

12528 de	1195 c	530 sera	341 ai	233 développement
8324 la	1188 je	528 doit	323 travail	231 économie
6211 l	1183 ne	527 aussi	310 entre	229 deux
5815 et	1127 par	509 ont	306 si	227 enfin
5217 les	1117 ce	494 français	297 économique	226 encore
4908 le	1074 sur	479 y	290 aujourd	226 temps
4631 à	985 qu	462 j	290 lui	222 ensemble
4435 des	908 france	453 etat	288 dont	221 vie
3832 d	855 s	447 sans	283 sociale	220 société
3051 est	838 aux	434 ou	282 on	219 depuis
2982 en	838 n	425 comme	280 seront	216 ceux
2799 que	816 nos	422 ces	278 monde	215 donc
2441 une	810 gouvernement	422 tout	278 république	210 toutes
2425 nous	803 avec	421 son	266 fait	209 soit
2273 qui	744 mais	413 avons	265 loi	208 droit
2142 un	711 elle	410 ses	265 où	208 sécurité
2060 pour	697 cette	409 même	264 contre	207 ainsi
2024 du	695 vous	406 été	263 leurs	206 elles
1977 dans	693 politique	400 faire	262 action	206 moyens
1809 il	667 se	390 ils	256 europe	203 cet
1410 au	651 être	386 faut	243 effort	202 autres
1393 notre	647 sont	375 entreprises	241 peut	202 cela
1368 plus	633 leur	362 emploi	236 nationale	199 mesures
1275 pas	603 pays	346 bien	235 avenir	197 jeunes
1214 a	533 tous	342 sa	235 président	195 croissance

12528 de	1195 c	530 sera	341 ai	233 développement	1147pre	ou	449con	savoir	312ver_sup
8324 la	1188 je	528 doit	323 travail	231 économie	1133pre	ces	449adj_dem	républicque	307nom
6211 l	1183 ne	527 aussi	310 entré		965con	falloir	444ver_sup	aujourd'hui	302adv_sup
5815 et	1127 par	509 ont	306 si		956ver_sup	vouloir	440ver_sup	année	297nom
5217 les	1117 ce	494 français	297 écon		948nr	comme	438con	on	294pro_per
4908 le	1074 sur	479 y	290 aujour		892pro_per	tout	433pro_ind	moyen	292nom
4631 à	985 qu	462 j	290 hui		885adj_pos	son	432adj_pos	européen	290adj
4435 des	908 france	453 état	288 dont		884ver_sup	public	432adj	dont	290pro_rel
3832 d	855 s	447 sans	283 social		881art_def	ses	431adj_pos	demande	277ver
3051 est	838 aux	434 ou	282 on		863pre	ils	411pro_per	économie	276nom
2982 en	838 n	425 comme	280 sero		849nom	économique	399adj	monde	276nom
2799 que	816 nos	422 ces	278 mon		811nom	premier	396adj	contre	273pre
2441 une	810 gouvernement	422 tout	278 répu		779con	national	390adj	société	270nom
2425 nous	803 avec	421 son	266 fait		722adj_dem	mettre	383ver	où	270pro_rel
2273 qui	744 mais	413 avons	265 loi		709pro_per	prendre	379ver	leurs	270adj_pos
2142 un	711 elle	410 ses	265 où		702adj	monsieur	379nom_sup	europe	270nr
2060 pour	697 cette	409 même	264 cont		695pro_per	bien	376nom_sup	mesure	266nom
2024 du	695 vous	406 été	263 leurs		694adj	travail	374nom	agir	266ver
1977 dans	693 politique	400 faire	262 actio		686pro_per	nouveau	374adj	an	263nom_sup
1809 il	667 se	390 ils	256 euro		666pro_per	sa	357adj_pos	président	258nom
1410 au	651 être	386 faut	243 effor		638nom	permettre	355ver	aller	257ver
1393 notre	647 sont	375 entreprises	241 peut		552adv_sup	effort	351nom	objectif	251nom
1368 plus	633 leur	362 emploi	236 natio		546pro_ind	même	347pro_ind	jeune	245adj
1275 pas	603 pays	346 bien	235 aven		543nom	action	328nom	engager	245ver
1214 a	533 tous	342 sa	235 prés		525nom	si	320con	projet	244nom
					517adj	loi	318nom	avenir	244nom
					503pro_per	entre	318pre	temps	243nom
					493pro_per	droit	318nom	service	241nom
					476nr	donner	313ver	réforme	241nom
					467pre	dire	313ver_sup	assurer	241ver

Tableau lexical



Parties (variable-s / segments de textes)

Lexique :

- Tokenization
- Reconnaissance
- Lemmatisation
- Statuts statistiques



- Nombre d'occurrences
- Présence / absence



Questions de recherche et résultats d'analyses: une typologie

- A. Mesurer et caractériser la diversité des expressions et les controverses en ligne
- B. Comparer plusieurs corpus, repérer convergences et divergences de représentations, opinions et/ou idéologies dans différents espaces
- C. Caractériser un processus ou l'évolution diachronique d'un discours



A. Mesurer et caractériser la diversité des expressions et les controverses en ligne

Analyse des commentaires de pétitions

change.org



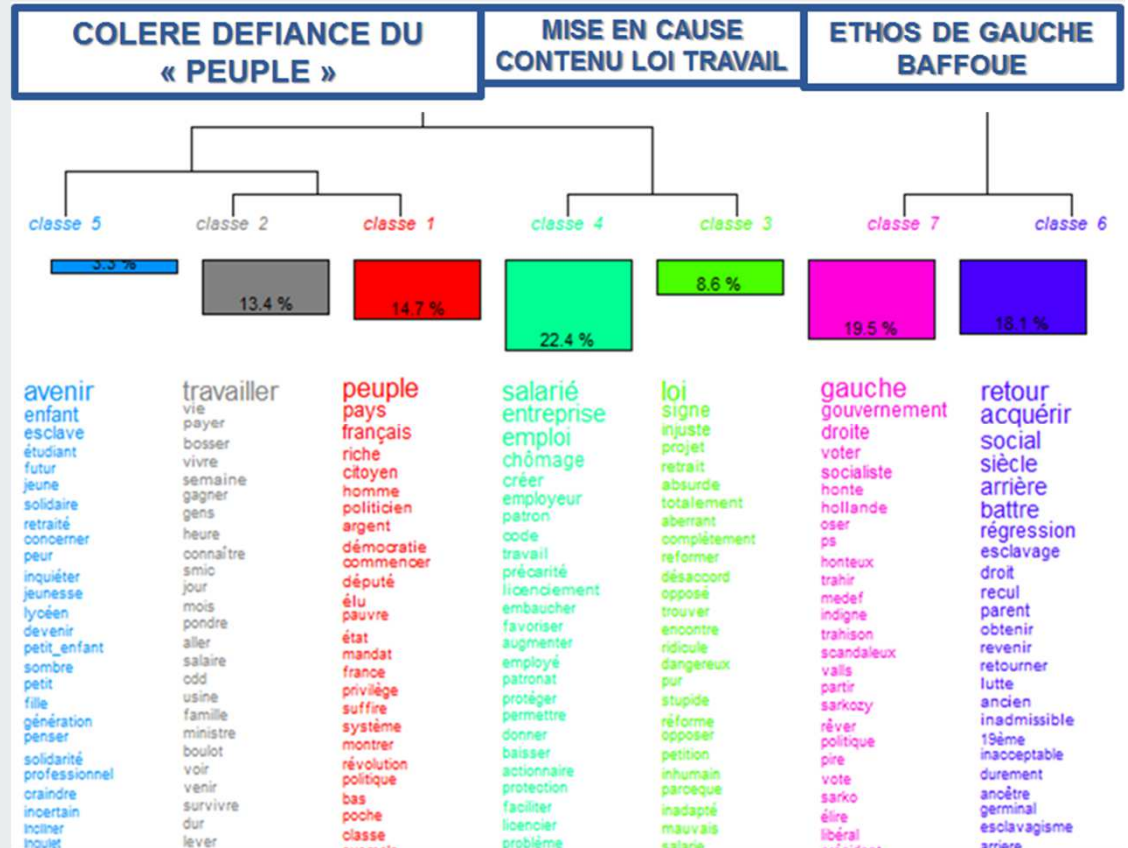
Bousquet, Marty, Smyrniotis (2017)

-286 836 commentaires du 19 fév au 31 mars
("Je signe parce que"...)

6 597 624 occurrences, 57797 formes lex.

→ **Diversité des expressions du mécontentement:**

- * **pathos:** colère du "peuple" contre les élites
- * **logos :** mise en cause argumentée de la loi Travail
- * **ethos:** choix politique mettant en danger l'histoire et les valeurs sociales de la France





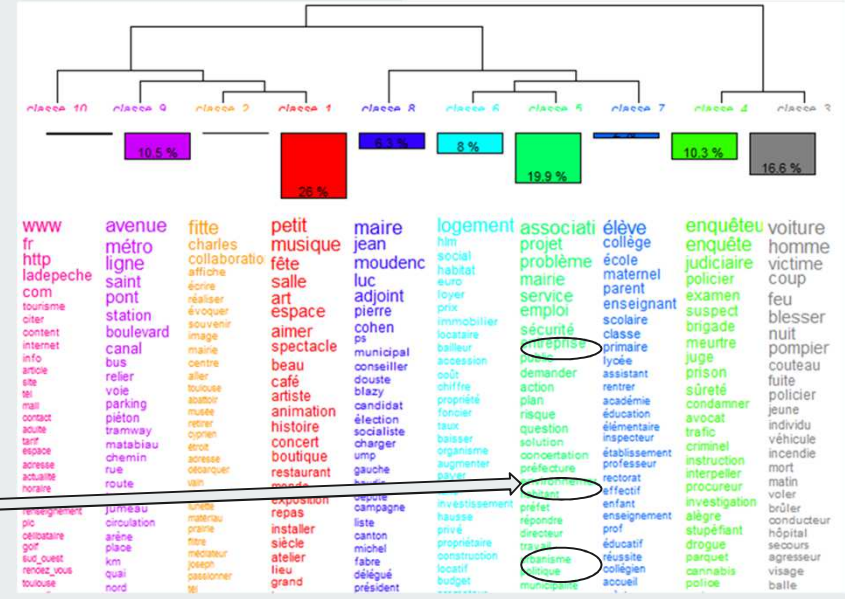
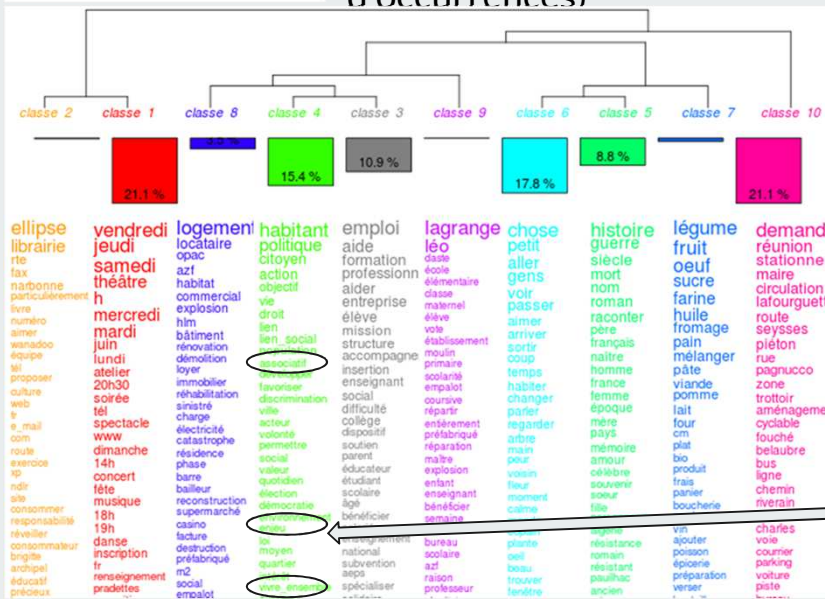
B. Comparer plusieurs corpus, repérer convergences et divergences

Différents médias

Natacha Souillard, 2018



Médias associatifs (10 journaux et bulletins de quartier : 2,6 millions occurrences) / PQR (4000 articles : 2,4 millions d'occurrences)



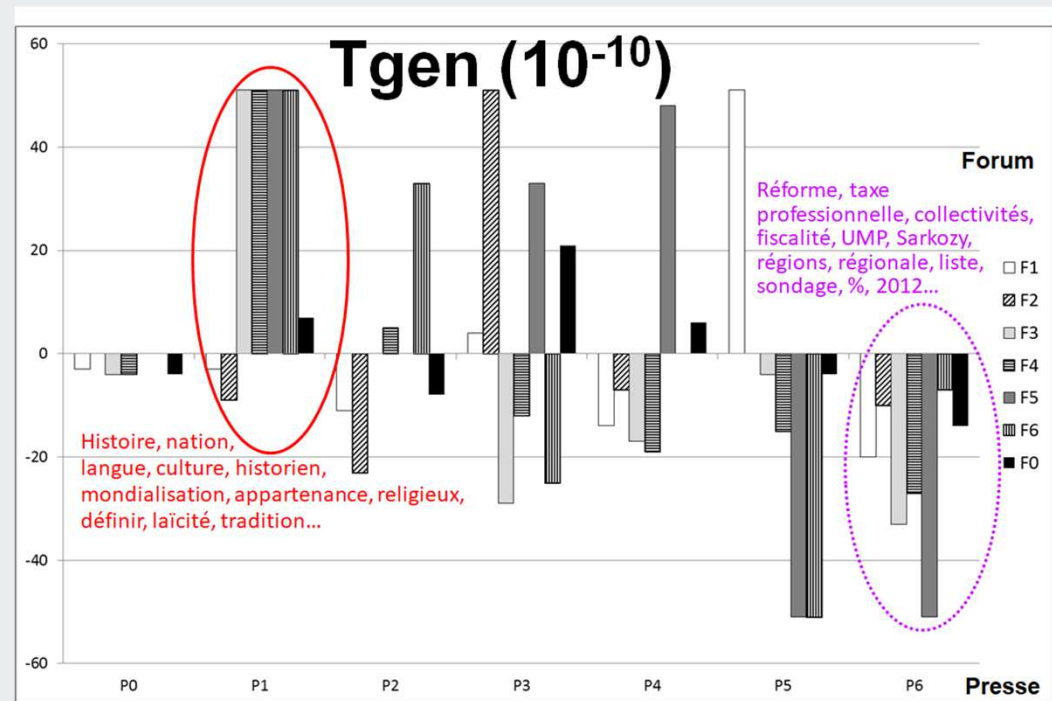
Le débat sur l'identité nationale

Forum : 18240 contributions

- 6 classes -> *Tgens*

Factiva : 1436 articles de presse

- 6 classes -> projection



(d'après Marty, Marchand & Ratinaud, *BMS*, 2012)



C. Caractériser un processus ou l'évolution diachronique d'un discours

La construction médiatique de la crise

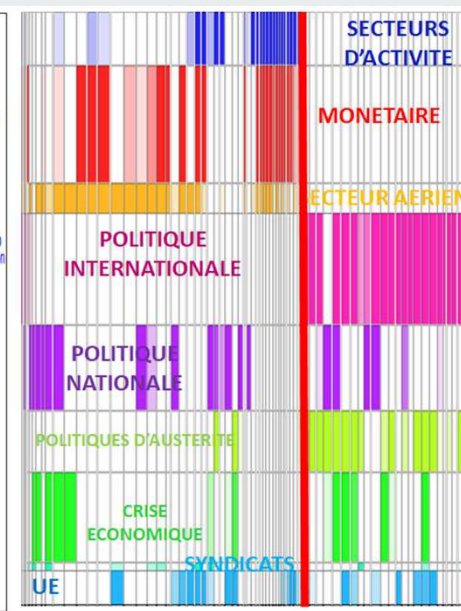
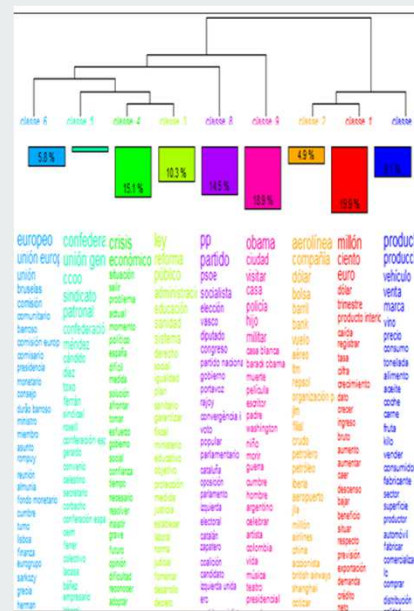
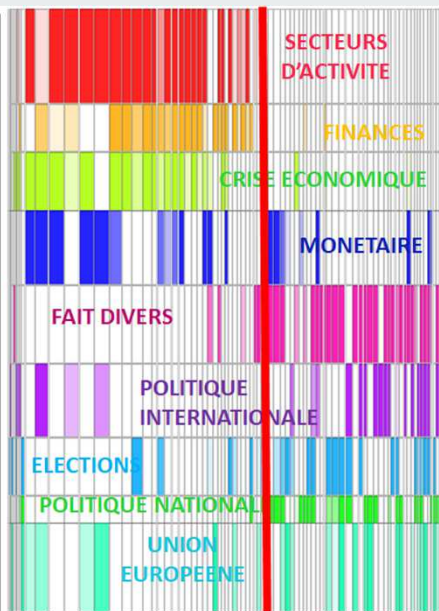
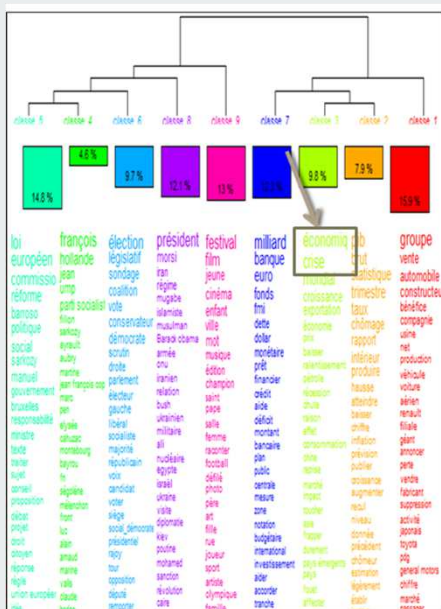
Mariola Moreno, 2018



17 850 textes, 7 962 140 occurrences, 66 579 formes

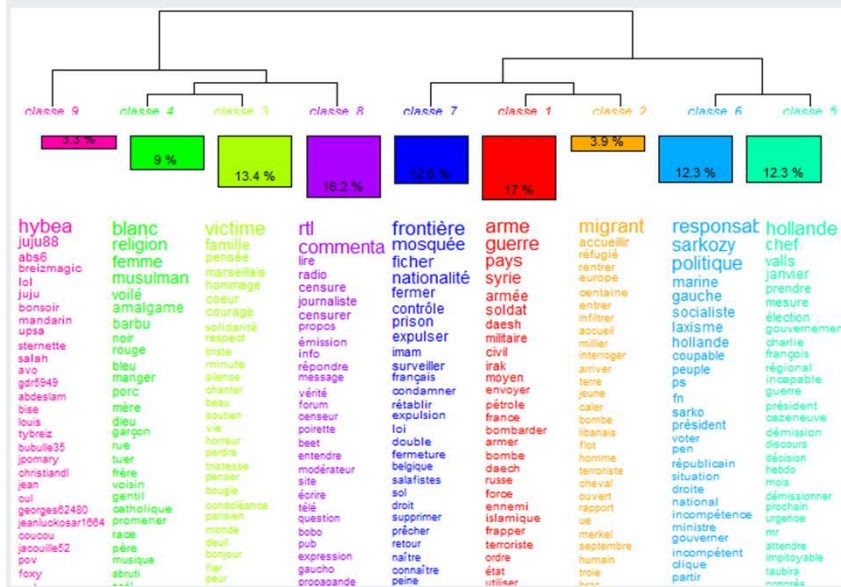


38 147 textes, 14 183 078 occurrences, 100 473 formes



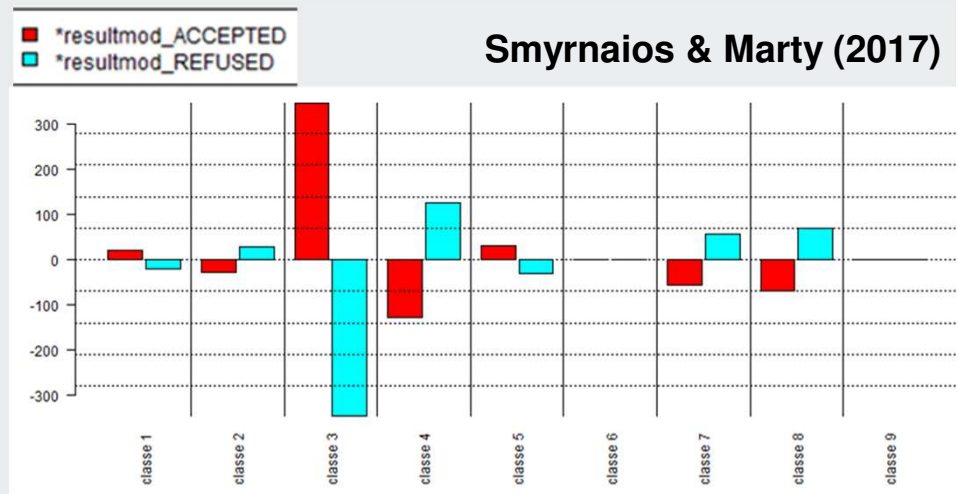
La modération de commentaires d'actualité

Commentaires apposés entre le 13 et le 16 novembre au bas d'articles traitant des attentats du 13 novembre, sur trois sites d'actualité, modérés en sous-traitance par la société Atchik. 29017 commentaires, 1 017 978 occurrences, 25943 formes lexicales



3 attitudes dans le processus de modération:

- refus immédiat (commentaires violents et/ou racistes)
- acceptation évidente (messages compassionnels et/ou descriptifs)
- examen minutieux et arbitrage (commentaires délicats, possiblement insultants ou diffamants)





Pour conclure:

Des outils statistiques pouvant être mobilisés pour différents contextes et corpus, et pour des questions de recherche différentes:

- La CHD comme première étape d'identification et de caractérisation d'une diversité d'expressions
- L'AFC pour identifier des distances/proximités, des polarisations et/ou des logiques structurantes
- La comparaison de profils de classe, facilitée par l'export de Types Généralisés pour mettre en regard des corpus différents, voire hétérogènes
- La mobilisation de variables indépendantes (ou méta-données) présidant à la partition pour identifier les liens significatifs entre éléments textuels et extra-textuels (locuteur ou source, temporalité, situation, caractérisation exogène du discours, etc.)



Pour conclure:

Statut de la statistique dans l'analyse de discours:

- administration d'une preuve statistique de l'organisation matérielle d'un discours (identité lexicale, spécificités, liens significatifs à certaines variables ou méta-données)
- dimension heuristique des indicateurs pour la reconstruction d'hypothèses interprétatives sur la signification des significativités, au regard des hypothèses et corpus mobilisés