ATHIS. L'historien, le texte et l'ordinateur. Les logiciels de traitement informatique du texte.

Alain Dallo Université Paris I. LAMOP

1. Des	logiciels qui ont une histoire. L'utilisation de corpus	4
1.1.	Plusieurs centres de recherche	4
1.2.	Des avancées récentes	4
1.3.	Corpus, logiciels, approches	4
2. L'ap	proche thématique	5
2.1.	La forme et sa distribution dans le corpus	5
2.1.1.	Index avec Hyperbase	6
2.1.2.	Dictionnaire des formes avec Lexico 3	7
2.1.3.	Dictionnaire des lemmes avec Weblex	8
2.1.4.		
2.1.5.	- · · · · · · · · · · · · · · · · · · ·	10
2.1.6.		
2.1.7.		
2.1.8.		
2.1.9.	7 1 6 1	
2.1.10	Rafale, carte des sections : Europe	15
2.2.	Les concordances	16
2.2.1.		17
2.2.2.	Concordances du lemme croire	18
2.2.3.	Concordances d'une liste de quelques formes avec Lexico 3	19
2.2.4.	Les concordances de deux lemmes éloignés avec Weblex. Utilisation de CQP = Corpus Query Processor	20
2.3.	Les segments répétés	22
2.3.1.		23
2.3.2.	Les segments répétés du corpus voeux. Weblex	24
2.4.	Les cooccurrences	25
2.4.1.		26
2.4.2.		
2.4.3.	L'environnement de la forme je	28
2.4.4.	Extrait de la liste de toutes les cooccurrences significatives	29
2.4.5.	Cooccurrences droites et gauches de France	30
2.4.6.	1 1 J	31
2.4.7	Résumé du graphique des cooccurrences de je	32

3. L'approche statistique ou lexicométrique	33	
3.1. Les caractéristiques d'un corpus	33	
3.1.1. L'étude des formes en fonction de leur fréquence	33	
3.1.2. Le nombre d'occurrences du corpus et de ses composantes	33	
3.1.3. Fréquence de nous	34	
3.1.4. Les mots outils dans le tableau lexical entier	35	
3.1.5. Les vrais et les faux hapax	36	
3.2. Les spécificités d'une partie	37	
3.2.1. J.C. Spécificités positives	38	
3.2.2. J.C. Spécificités négatives	39	
3.2.3. J.C. Spécificités phrastiques	40	
3.3. Rapprochement et éloignement des textes d'un corpus	41	
3.3.1. L'analyse factorielle du dictionnaire		
3.3.2. AFC sur une liste de mots	43	
3.3.3. Le classement automatique hiérarchisé	44	
4. Une autre approche en gestation	45	
4.1. Lemmes et catégories grammaticales	45	
4.1.1. formes et lemmes		
4.1.2. formes et catégories	47	
4.1.3. concordance dnav n=France	48	
4.1.4. concordance (N, V, A) nom = France, verbe = être,	49	
5. Quelques liens	50	
5.1. Logiciels	50	
5.2. Corpus	50	

1. Des logiciels qui ont une histoire. L'utilisation de corpus.

1.1. Plusieurs centres de recherche

- Charles Muller
- Étienne Brunet. Université de Nice
- ENS de Saint-Cloud
- Maurice Gross, Paris VII. Max Silberztein

1.2. Des avancées récentes

- Corpus étiquetés
- Préparation des corpus et balisage morpho-syntaxique
- Les moteurs de recherche, la fouille de textes et l'indexation fine de documents

1.3. Corpus, logiciels, approches

Partir du mot pour aller jusqu'au positionnement des textes d'un corpus, en passant par des études thématiques

- Un corpus : les voeux des présidents de la V ème République.
 - ° Corpus mis en place sur Weblex dans le cadre de la thèse de Jean-Marc Leblanc
 - ° Corpus intégré dans les autres logiciels à partir des textes fournis par Damon Mayaffre
- ♦ Les logiciels : Hyperbase, Lexico3, Weblex et Nooi
- ♦ L'approche thématique ou documentaire
- ♦ L'approche statistique

2. L'approche thématique

De la forme isolée à son environnement

- L'étude du dictionnaire
- Les concordances d'une forme
- Les segments
- Les cooccurrences

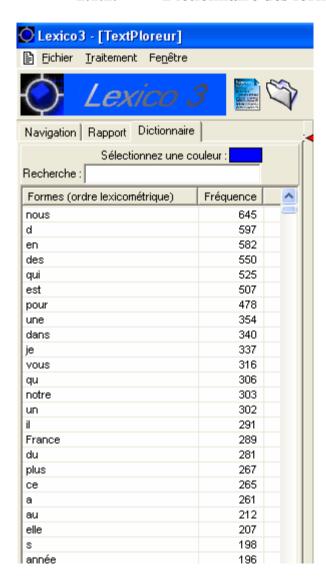
2.1. La forme et sa distribution dans le corpus

- Un index avec Hyperbase
- Un dictionnaire des formes avec Lexico 3
- Dictionnaire des lemmes avec avec Weblex
- Le dictionnaire des lemmes repérés : Nooj. Images 1 et 2.
- Calcul des rafales des voeux dans Weblex et graphique de la répartition de je
- Rafales et distribution d'une forme avec Hyperbase. Images : 1 et 2
- Rafales et distribution d'une forme avec Lexico3

2.1.1. Index avec Hyperbase



2.1.2. Dictionnaire des formes avec Lexico 3



2.1.3. Dictionnaire des lemmes avec Weblex

fréq.

4

4

5 12

3

9

8

4

Vocabulaire «élagué» du (

p3s de fréquence supérieure ou égale à 3.

Liste	alphabétique
-------	--------------

ord p3

25

26

28

29

30

31

32

Liban

Nouvel

Paris

Pologne

Duggio

Noël

Moyen-Orient 27 Nations-Unies

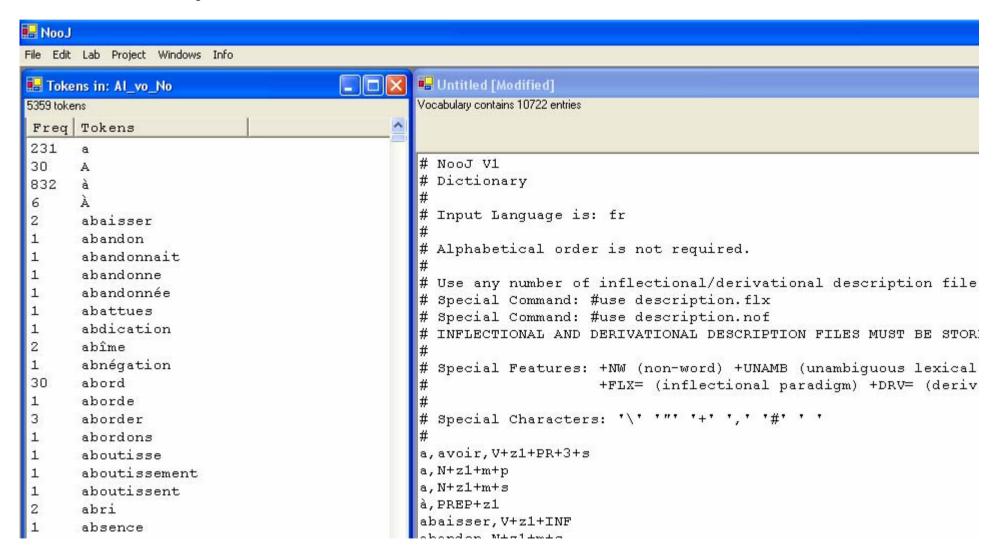
Outre-Mer

1	*	18
2	Afrique	14
3	Airbus	4
4	Algérie	21
5	Allemagne	12
6	Amérique	15
7	Angleterre	3
8	Asie	6
9	Berlin	3
10	Bosnie	3
11	Canada	3
12	Chine	4
13	Douze	3
14	Europe	99
15	France	300
16	Gaulle	3
17	Gorbatchev	3
18	Grande-Bretagne	3
19	Guerre	3
20	Irak	4
21	Israël	3
22	Italie	4
23	Japon	4
24	Koweït	5

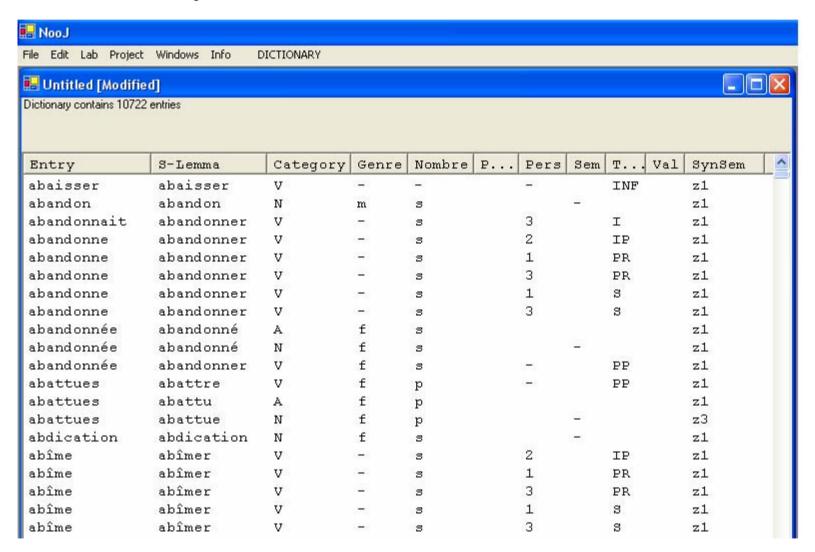
Liste hiérarchique

р3	fréq.
nous	654
je	426
France	300
français	219
année	202
faire	183
pouvoir	149
devoir	125
pays	122
on	113
monde	102
Europe	99
vouloir	92
voeu	88
bon	87
grand	85
nouveau	85
dire	82
peuple	82
an	80
souhaiter	79
savoir	77
aller	75
falloir	68
vivre	68
état	67
cher	66
vie	66
paix	64
compatriote	63
politique	61
République	60
comic!	50

2.1.4. Nooj: liste des formes et dictionnaire



2.1.5. Nooj: Liste des formes et dictionnaire

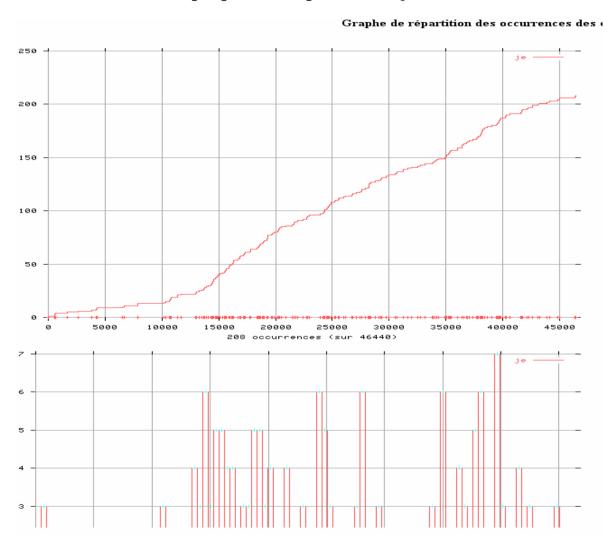


2.1.6. Calcul des rafales dans weblex

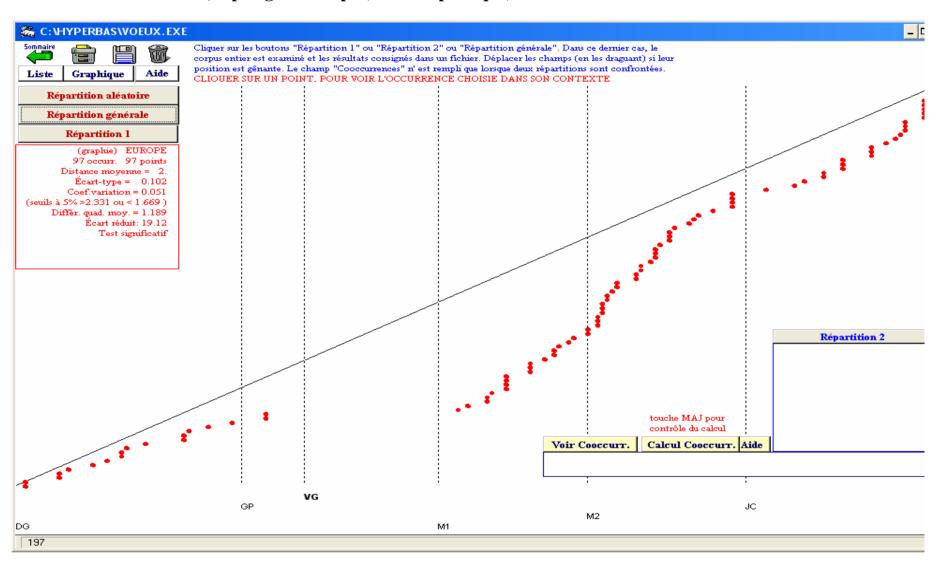
Répartition du vocal

				Repartition du vocai
	Z décroissant			
	forme	z	fp	
1	compatriote	20.8582	63	
2	euro	20.145	13	
3	je	17.1565	376	
4	Europe	14.4225	91	
5	Algérie	14.1149	19	
6	nier	13.8835	10	
7	peur	13.8149	8	
8	cher	13.0681	65	
9	scientifique	12.8823	10	
10	nous	12.0972	483	
11	but	11.8479	10	
12	progrès	11.6406	51	
13	ami	11.5533	11	
14	développement	10.6219	22	
15	difficultés	10.6133	19	
16	soldat	9.92321	10	
17	vive	9.28211	53	
18	apporter	9.22639	23	
19	hausse	8.88236	6	
20	au-dehors	8.38362	7	
21	*	8.18859	10	
22	bonheur	8.11869	20	
23	valeurs	8.08423	11	
24	démocratie	8.01422	16	
25	fête	7.25487	7	

2.1.7. Graphique de la répartition de je



2.1.8. Rafale, topologie : Europe (Vème République)



2.1.9. Rafale, topologie: Europe chez Mitterand

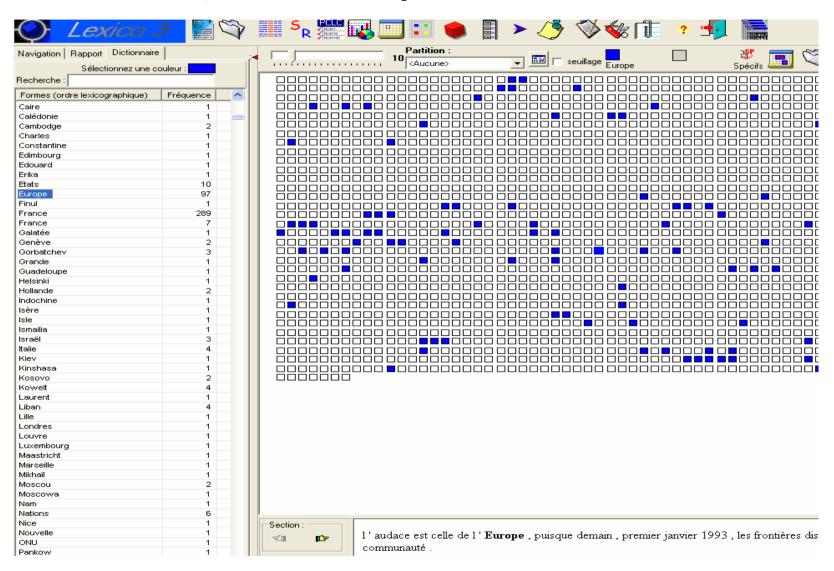


Depuis quelques mois , nous avons dans les yeux les images terribles des combats que se livrent les peuples de Yougoslavie , hier encore associés sous un même drapeau . Comment arrêter cette guerre ?

La France soutient les efforts de négociation et d'arbitrage de la Communauté et des Nations - Unies . Elle reconnaît le principe de l'autodétermination . Mais il lui paraît urgent que soient mises en place des structure inter - européennes , où le droit à l'indépendance ne se confondra pas avec l'anarchie des tribus d'autrefois . Ce sera , je le pense , l'un des enjeux majeurs de 1992 .

Car l' inquiétude gagne l' Europe de l' Est où l' on redoute la contagion . Comment cette inquiétude nou épargnerait - elle , nous qui , à l' Ouest , avons pourtant la chance de vivre en paix et d' avoir dépassé nos propre divisions ? Raison de plus de se réjouir des récents accords de Maastricht . Une monnaie commune , l' amorce d' une diplomatie , d' une défense et d' une armée communes à l' Europe des Douze , une charte sociale , l' exemple de stabilité offert aux peuples qui se déchirent , bientôt 350 à 360 millions d' Européens solidaires sur la scène de monde - en attendant les autres - , bref , l' Europe qui se fait , voilà un grand dessein , capable d' enthousiasmer , de rassembler et de justifier l' espérance .

2.1.10. Rafale, carte des sections : Europe

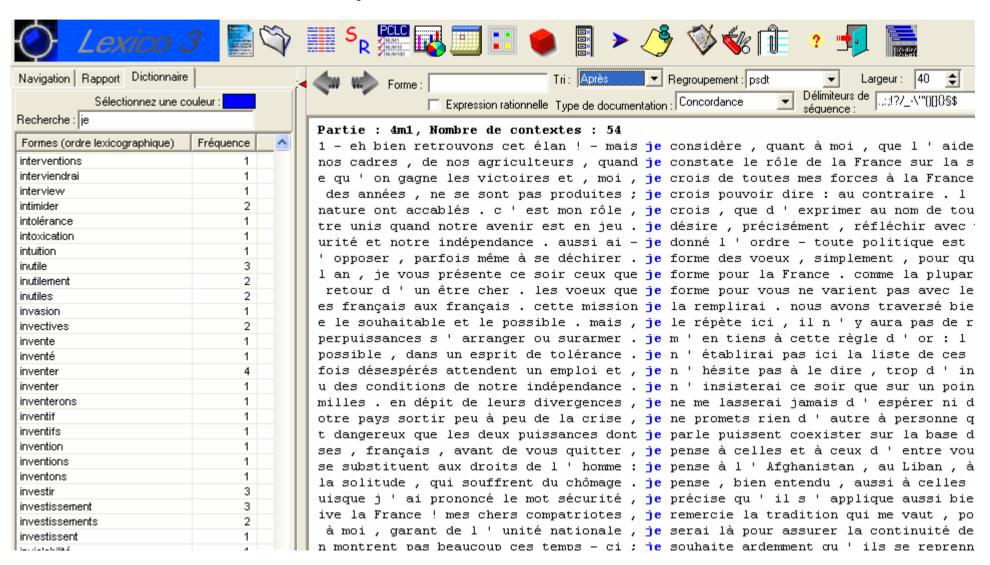


2.2. Les concordances

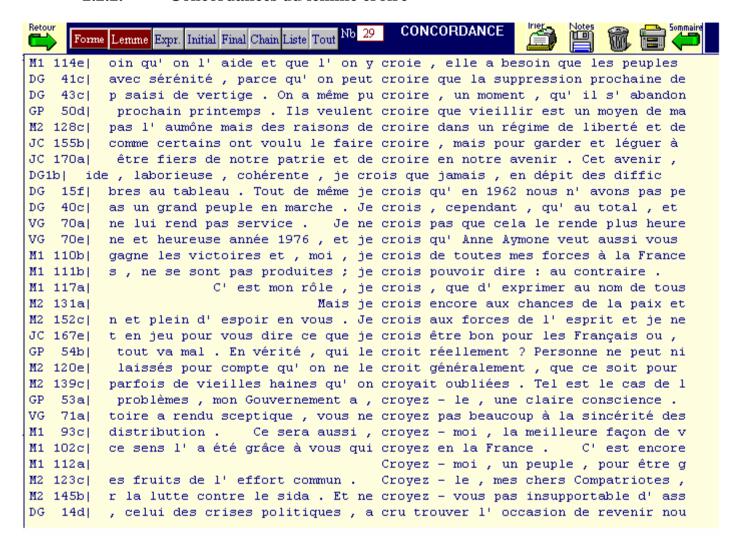
L'environnement proche d'une forme

- Les concordances d'une forme : Lexico 3
- Les concordances d'un lemme avec Hyperbase
- Les concordances d'une liste de formes
- Les concordances de deux lemmes éloignés avec Weblex. Utilisation de CQP = Corpus Query Processor

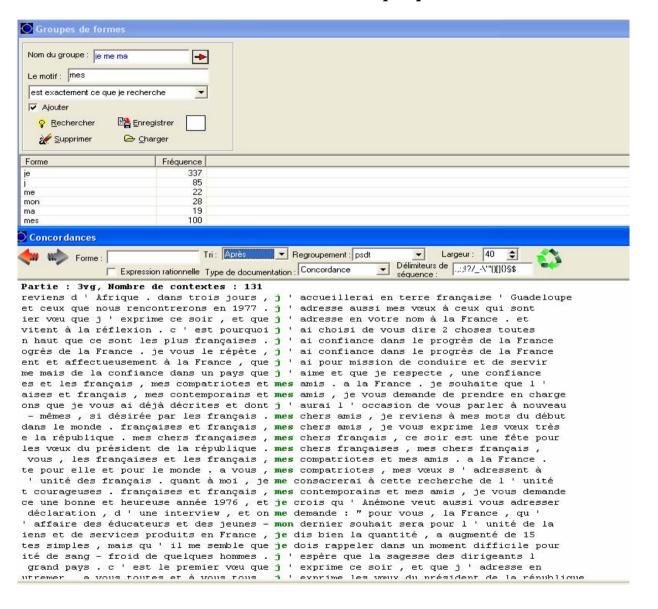
2.2.1. Concordances de la forme je chez Mitterand



2.2.2. Concordances du lemme croire



2.2.3. Concordances d'une liste de quelques formes avec Lexico 3



2.2.4. Les concordances de deux lemmes éloignés avec Weblex. Utilisation de CQP = Corpus Query Processor

Chercher dans la page :]

Il y a 7 occurrences de "France"[]*"enfant" dans le corpus voeux2-author-Degaulle (propr

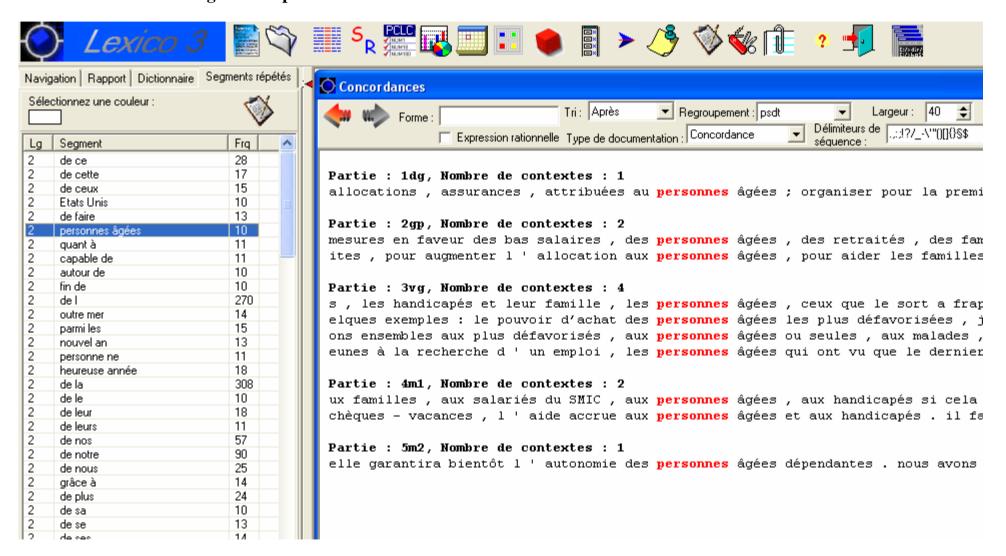
- 1	<u>Degaulle,</u> 1960	ur et à son énergie . Car l' Algérie a besoin de la communauté française , et la	France , pour son oeuvre , a besoin d' elle en Algérie . Bien entendu , et , quoi qu' il arrive , la enfants
- 1	0egaulle <u>.</u> 1963	ustes envers nous-mêmes . Car le bilan est positif . Pendant ces douze mois , la	France a continué de monter . En 1963 , notre population s' est augmentée de presque six cen est né environ neuf cent mille bébés . La proportion croissante des jeunes devant progressiv processus , on peut penser que , parmi les enfants
	<u>Degaulle,</u> 1 <u>963</u>	ns, j' ai l'honneur et la charge de parler en notre nom à tous, j' offre à la	France , cette fois encore , les voeux très ardents et très confiants de ses enf
- 1	Degaulle, 1964		France . Certes , la vie est la vie , autrement dit un combat , pour une nation comme pour un ho , toujours et partout , des difficultés à vaincre , des efforts à déployer , des peines à supporter , de dignité , de justice , de fraternité . Mais , tous ensemble , nous sommes un peuple , dont l' év les enfants
- 1	Degaulle <u>.</u> 1965	ns nos têtes , dans nos coeurs et dans nos mains tout ce qu' il faut pour que la	la France parcoure une étane décisive de son progrés a pour ou'elle apporte plus de justice enco-
- 1	Degaulle <u>.</u> 1965	eurs voeux pour 1966 et que , tous ensemble nous souhaitons une bonne année à la	France . Vive la République ! Vive la France . Françaises , Français , Je vous souhaite une bon nom de nous tous , mes voeux sont l'expression de ceux que la France forme pour chacu

2.3. Les segments répétés

Les successions de formes contiguës du texte

- Avec Lexico 3
- Avec Weblex

2.3.1. Les segments répétés avec Lexico 3



2.3.2. Les segments répétés du corpus voeux. Weblex

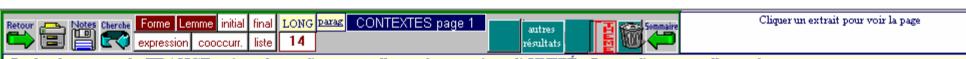
f segment	l
8 chacune et à chacun d'entre vou	<u>18</u> 7
5 à chacune et à chacun d'entre	7
5 à chacune et à chacun de vous	7
3 C' est au nom de l' emploi	7
3 au nom de l'emploi que nous	7
3 est au nom de l'emploi que	7
3 À chacune et à chacun d'entre	7
8 chacune et à chacun d'entre	6
8 et à chacun d'entre vous	б
6 à chacune et à chacun de	6
5 chacune et à chacun de vous	б
5 à chacune et à chacun d'	6
4 Nous ne sommes pas les plus	б
4 À chacune et à chacun d'	6
3 C' est au nom de l'	б
3 II n' y a pas d'	6
3 Personne ne peut nier que la	б
3 au nom de l'emploi que	6
3 ceux qui ont la chance d'	б
5 4 40 415	-

2.4. Les cooccurrences

Présence de deux formes non contiguës

- Les cooccurrences de deux formes dans hyperbase
- Utilisation de la fonction thème d'hperbase pour les découvrir. Images : 1 et 2.
- Le calcul des cooccurrences dans weblex
- Cooccurrences gauches et droites
- Graphique
- Résumé du graphique

2.4.1. Les cooccurrences : Unité France



Le bonheur pour la FRANCE, c'est la confiance en elle - même et c'est l'UNITÉ. La confiance en elle - même, parce que nous traverson une époque difficile, qui est une époque d'évolution et d'adaptation dans le monde, à la recherche d'un nouvel équilibre. La FRANCE doi avoir confiance en elle - même, parce qu'elle est capable de surmonter ces difficultés. Quand nous pensons a ce qui s'est passe en FRANCI depuis le début de ce siècle, et dont beaucoup d'entre vous se souviennent, quand nous pensons a la première guerre qui a décimé la plupar des familles françaises, quand nous pensons au deuxième conflit dans lequel notre pays a été totalement occupe et son économie désorganise et lorsque nous constatons que nous avons été capables de surmonter ces difficultés, la FRANCE peut avoir confiance en elle - même pou surmonter les difficultés du présent.

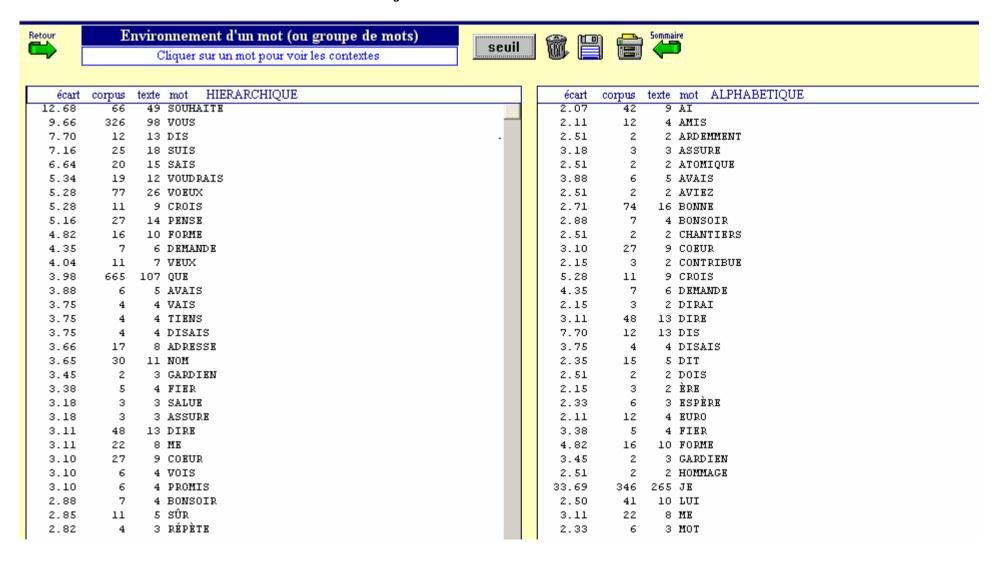
mon dernier souhait sera pour l' UNITÉ de la FRANCE . La FRANCE est au total , malgré d'inutiles querelles , plus unie qu'il y a un an . Or pouvait craindre qu'elle ne sorte déchirée du grave débat politique de mars dernier `élections législatives` . La haine et la rancune pouvaient nouveau diviser les français . Mais le sens de l' UNITÉ ce sens venu du fond des ages et qui est le certificat de naissance des nations le sens d' UNITÉ l' a emporte grâce au _concours , je tiens a le dire , des uns et des autres .

En 1980 , la FRANCE a besoin du sérieux , du courage , des facultés d'adaptation , de la générosité et de l'UNITÉ des Français .

2.4.2. Les contextes courts de je



2.4.3. L'environnement de la forme je



2.4.4. Extrait de la liste de toutes les cooccurrences significatives

Cooccurrences du corpus "voeux2" (forme)

Seuils: f 3, cf 3, p 5.0E-2, d_m 1000.0

AB	fa frecf p do	
	-R -D P	
Chers Compatriotes	34 34 34 2e-75 0.0	
Françaises Français	38 124 34 8e-41 0.7	
chers compatriotes	21 28 18 1e-34 0.0	
Vive République	55 60 29 3e-34 1.0	
bonne année	51 199 37 4e-30 0.7	
Nouvel An	12 16 12 1e-28 0.0	
Je souhaite	136 65 34 2e-26 0.2	
personnes âgées	12 10 10 1e-25 0.0	
pouvoir achat	18 9 9 2e-20 1.0	
intérêts particuliers	15 10 9 2e-20 0.0	
Bonne année	24 199 21 2e-20 0.4	
bonne heureuse	51 22 14 6e-19 1.6	
heureuse année	22 199 19 3e-18 1.6	
souhaite bonne	65 51 18 3e-16 7.0	
je souhaite	208 65 30 5e-16 0.2	
peut nier	41 10 9 1e-15 0.0	
cent mille	12 10 7 3e-15 0.4	
départements territoires	6 6 5 1e-13 1.4	

2.4.5. Cooccurrences droites et gauches de France

Lexicogramme du pôle "France" dans le corpus voeux2-author-Degaulle (p3)

```
Seuils: f 3, cf 3, p 5.0E-2, d<sub>m</sub> 1000.0
                     France
                      (70)
   cooccurrents gauches cooccurrents droits
            fcfpd<sub>m</sub>
                                  fcf p d<sub>m</sub>
          18 9 1e-04 1.0 enfant 8 5 1e-03 8.6
vive
souhaiter 12 6 2e-03 4.5
          10 5 5e-03 13.2
nom
adresser 4 3 8e-03 14.7
<u>année</u> <u>47 11</u> 2e-02 6.5
         10 4 3e-02 13.0
voeu
maintenant 6 3 3e-02 16.3
           6 3 3e-02 16.7
terre
```

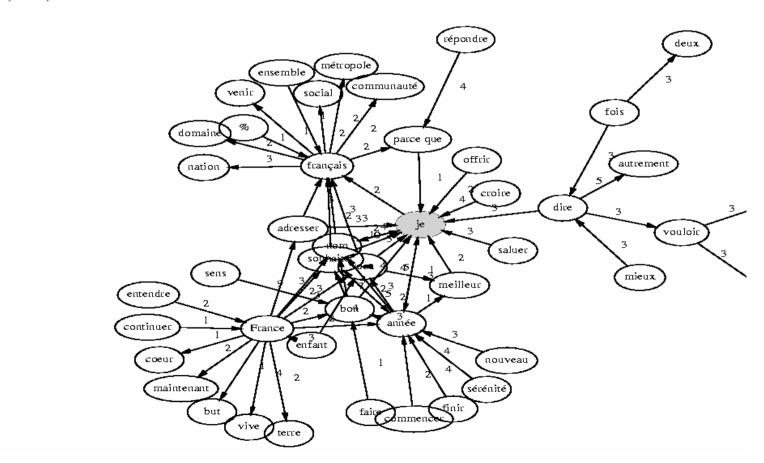
2.4.6. Graphique des cooccurrences de je

Lexicogramme récursif autour du pôle je dans le corpus voeux2-author-Degaulle

Seuils: p 2e-01, r 3, f 3, d_m 1000.0, pl 3

Synthèse

0 noeuds (0 arcs)



2.4.7. Résumé du graphique des cooccurrences de je

Calcul en cours (abandon)...

Synthèse des lexicogrammes récursifs du corpus voeux2-author-Degaulle

Seuils: \mathbf{p} 5e-02, \mathbf{r} 3, \mathbf{f} 3, $\mathbf{d}_{\mathbf{m}}$ 1000.0, \mathbf{pl} 3

Tri décroissant par le nombre de noeuds par lexicogramme.

- 28 (nous avancer prospérité voie développement Amérique latin disposer lui-même part prendre propre évolution aide état Europe détente union coopération puis grand nombre peuple Afrique)
- 12 (communauté France finir nouveau sérénité bon je nom souhaiter année commencer Algérie)
- 11 (enfant meilleur voeu métropole nation catégories domaine économique social français adresser)
- 5 (d abord expansion niveau vie agir)
- 3 (commerce industrie agriculture)

3. L'approche statistique ou lexicométrique

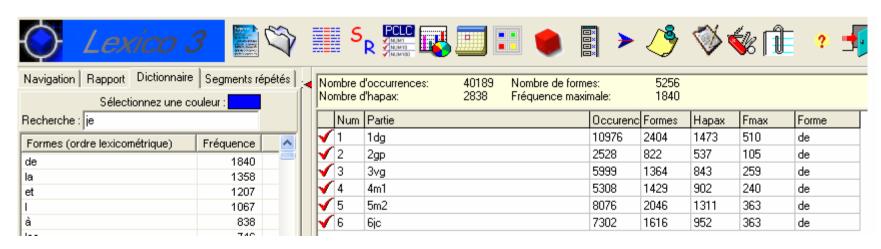
De l'étude des fréquences des formes au positionnement des textes

- Les caractéristiques d'un corpus
- Les spécificités d'une partie de celui-ci
- Rapprochement et éloignement des textes le composant

3.1. Les caractéristiques d'un corpus

- 3.1.1. L'étude des formes en fonction de leur fréquence
 - Le nombre d'occurrences
 - Les fréquences absolues et relatives
 - Les mots outils
 - Les hapax

3.1.2. Le nombre d'occurrences du corpus et de ses composantes

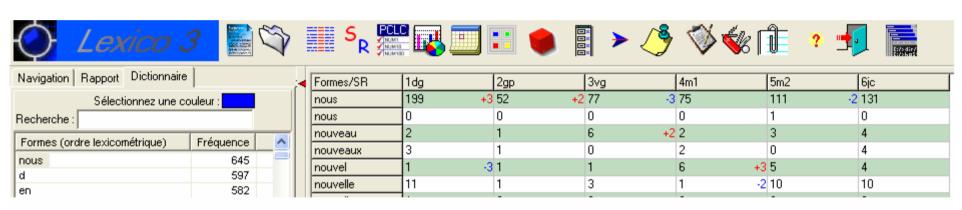


3.1.3. Fréquence de nous

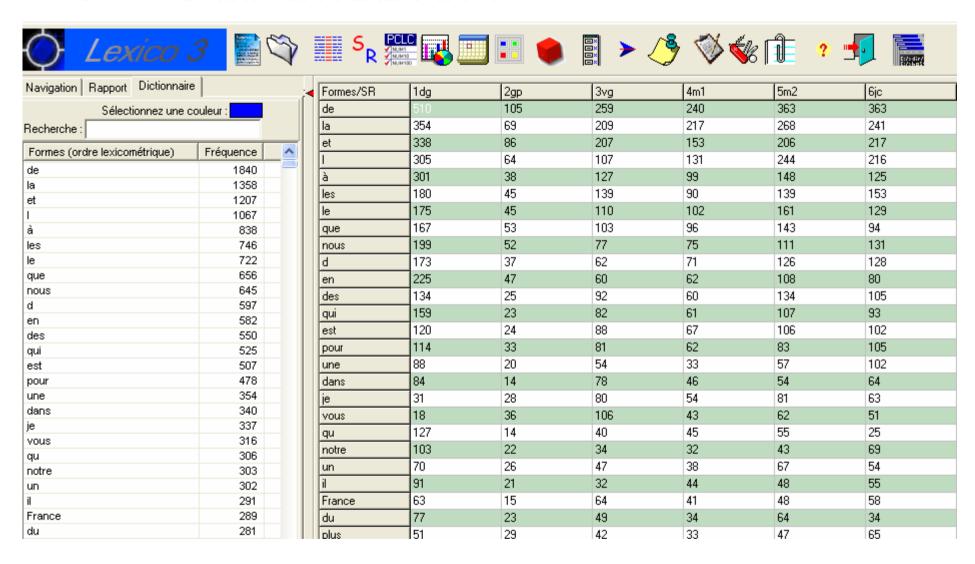
Nb. d'occurrences du corpus : 40189

Fréquence de nous : 645

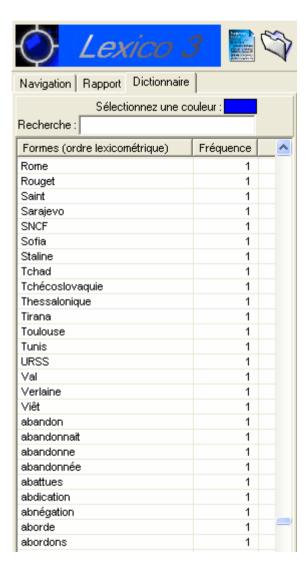
Fréquence relative de nous : 645 / 40189 = 0.016Cadence : 40189 / 645 = 62



3.1.4. Les mots outils dans le tableau lexical entier



3.1.5. Les vrais et les faux hapax

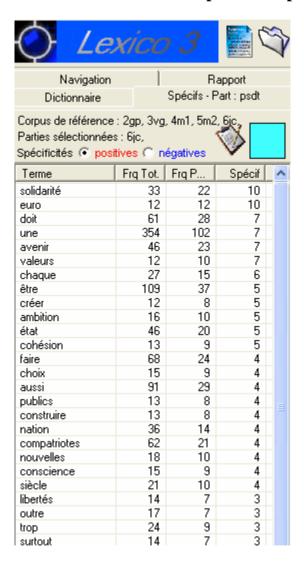


3.2. Les spécificités d'une partie

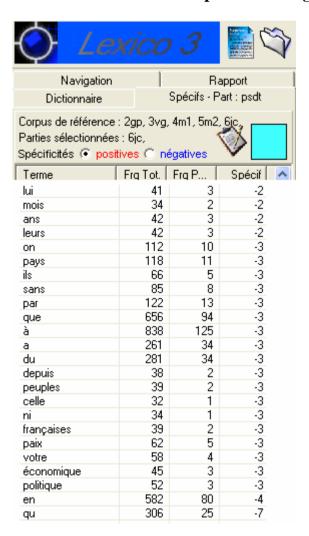
Les formes caractérisant une partie d'un corpus

- La loi hypergéométrique
- J.C. Spécificités positives et négatives
- J.C. Spécificités phrastiques

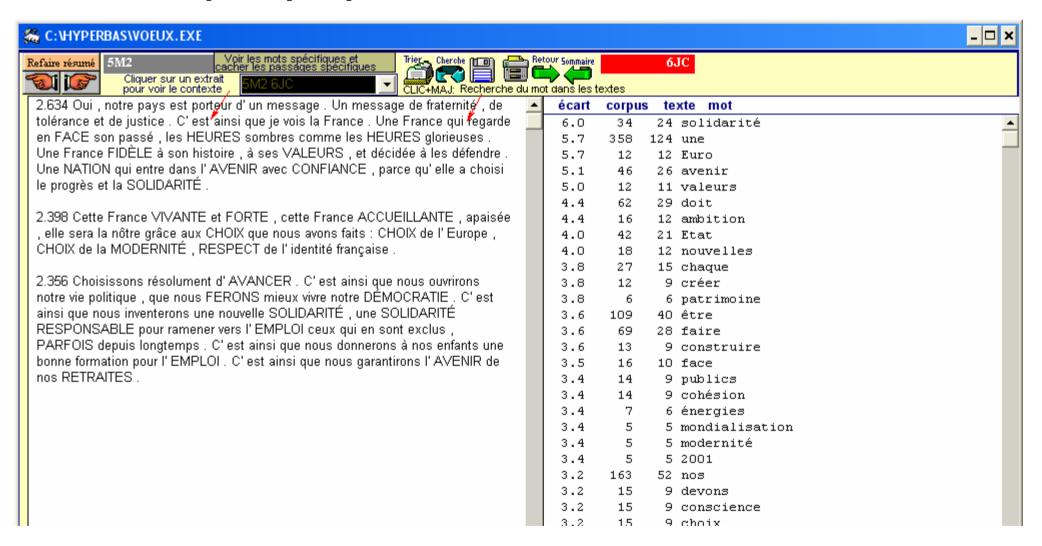
3.2.1. J.C. Spécificités positives



3.2.2. J.C. Spécificités négatives



3.2.3. J.C. Spécificités phrastiques

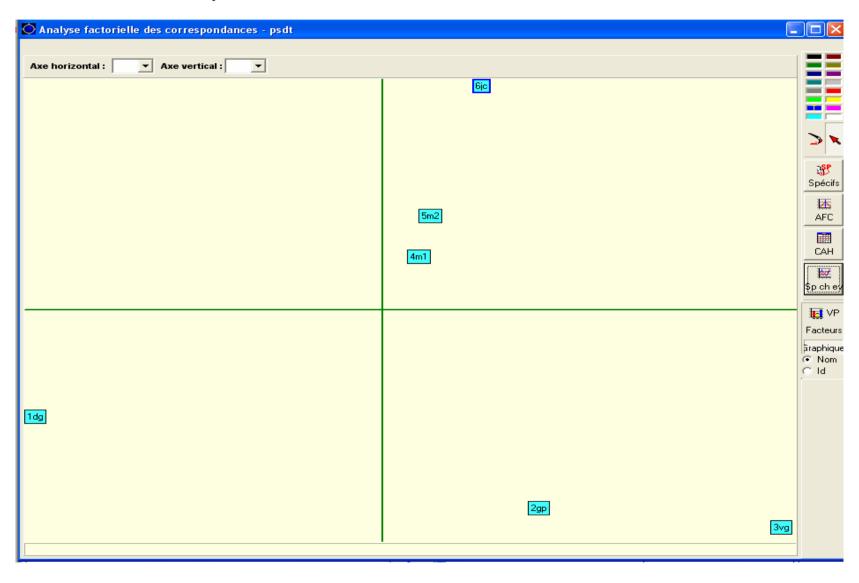


3.3. Rapprochement et éloignement des textes d'un corpus

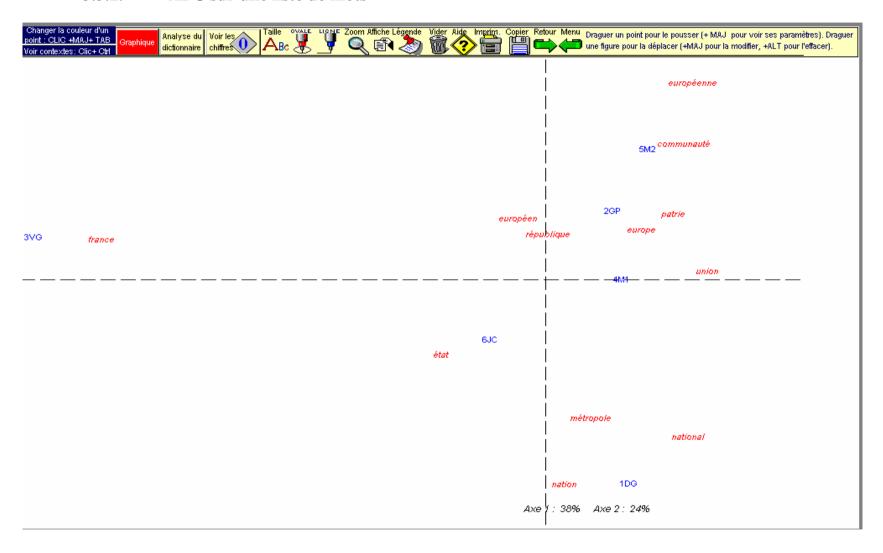
Oppositions, rapprochements et distances entre les textes d'un corpus

- L'analyse factorielle du dictionnaire
- AFC sur une liste de mots
- Le classement automatique hiérarchisé

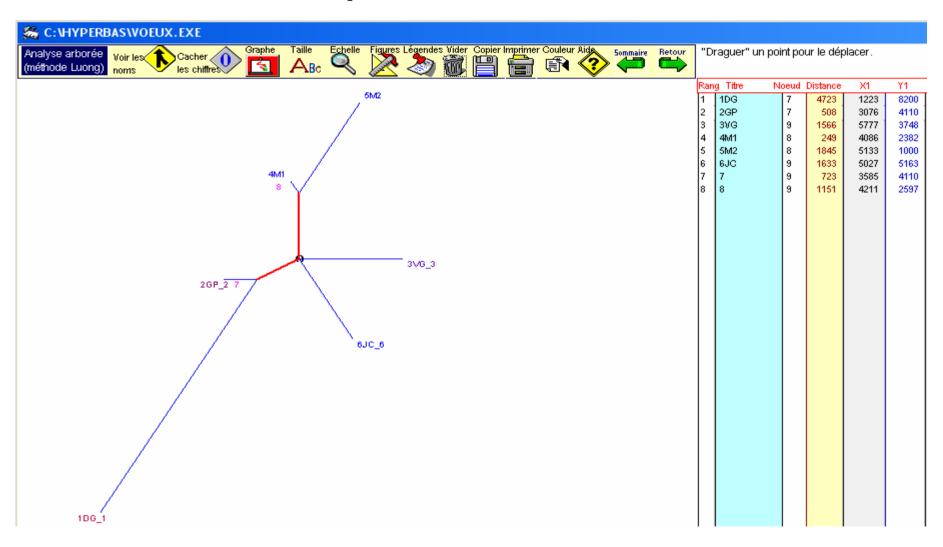
3.3.1. L'analyse factorielle du dictionnaire



3.3.2. AFC sur une liste de mots



3.3.3. Le classement automatique hiérarchisé



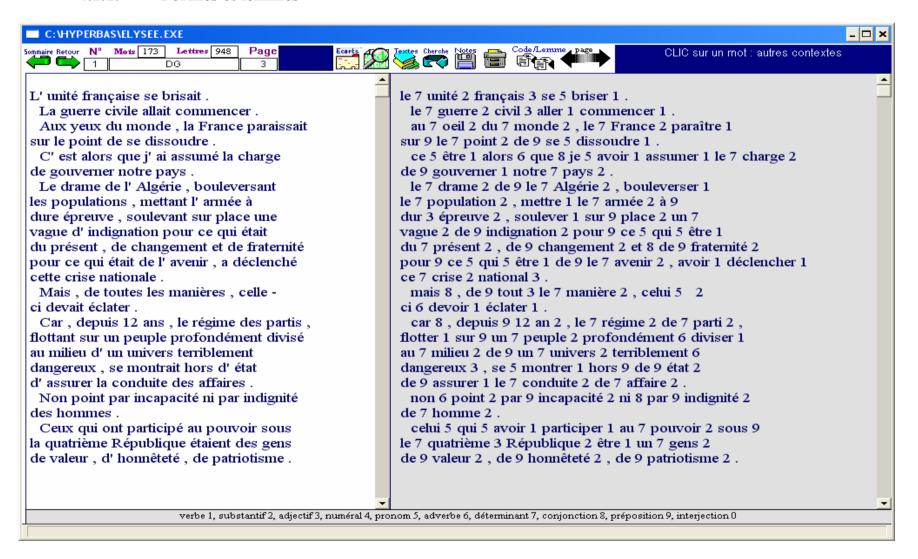
4. Une autre approche en gestation

- L'intégration de corpus et de textes formatés différemment dont ceux en xml (TEI)
- Les traitements préalables des textes : treetagger, cordial
- Un travail plus tourné sur des formes lemmatisées
- Des analyses plus fines des formes et des segments
- Un ou des logiciels exploitant un fond commun de petits outils

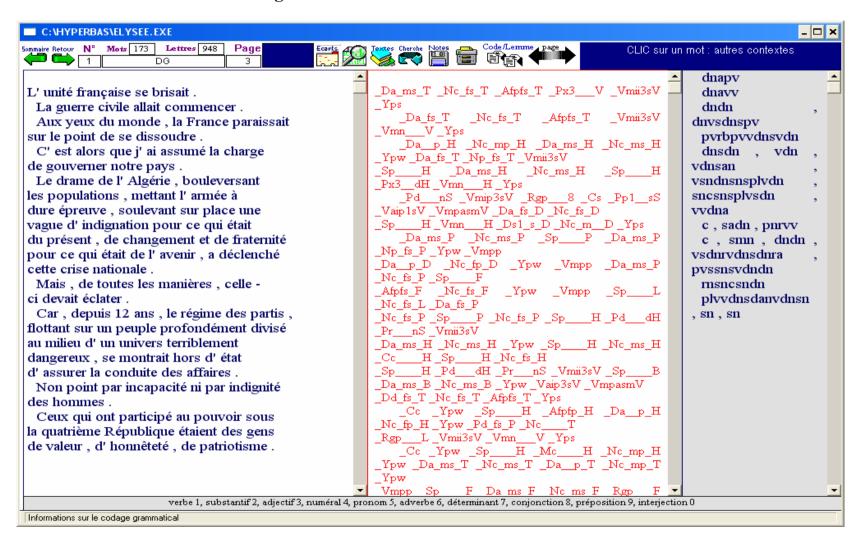
4.1. Lemmes et catégories grammaticales

- Hy : Formes non lemmatisées et lemmes
- Hy: Formes et catégories grammaticales
- Hy: concordance D N A V (det. nom. adj. ver.) nom = France
- Nooj: concordance (N, V, A) nom = France, verbe = être,

4.1.1. Formes et lemmes



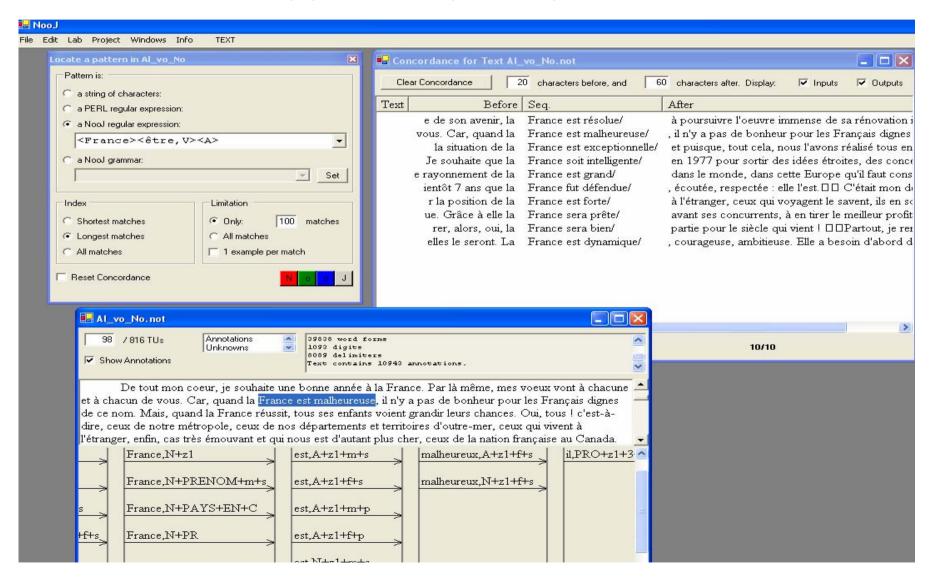
4.1.2. Formes et catégories



4.1.3. Concordance dnay n=France



4.1.4. Concordance (N, V, A) nom = France, verbe = être,



5. Quelques liens

- Logiciels
- Corpus

5.1. Logiciels

- Hyperbase: Etienne Brunet, UFR Lettres, 98 bd Herriot, 06204 Nice Fax: (33) 04 93 37 54 45 Courriel: brunet@unice.fr
- <u>Lexico 3</u> http://www.cavi.univ-paris3.fr/Ilpga/ilpga/tal/lexicoWWW/
- Weblex http://weblex.ens-lsh.fr/doc/weblex/index.html#sommaire
- Nooj http://perso.orange.fr/rosavram/pages/noojpag.html

5.2. Corpus

- <u>Damon Mayaffre : Politext</u> http://www.unice.fr/ILF-CNRS/politext/
- <u>La documentation Française</u> http://discours-publics.vie-publique.fr/rechlogos/servlet/RechServlet?_page=ACCUEIL&_type=NEW
- <u>Textopol</u> http://textopol.free.fr/
- <u>Jean-Marc Leblanc</u> http://www.univ-paris12.fr/www/labos/ceditec/leblanc.html/