



HAL
open science

Plaidoyer pour la statistique linguistique

Étienne Brunet

► **To cite this version:**

Étienne Brunet. Plaidoyer pour la statistique linguistique. Céline Poudat. Ce qui compte. Ecrits choisis tome II, 2, Champion, pp.311-329, 2011, 978-2-7453-2225-8. hal-01580838

HAL Id: hal-01580838

<https://hal.science/hal-01580838>

Submitted on 2 Sep 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Plaidoyer pour la statistique linguistique

Etienne Brunet

Charles Muller, le maître de la statistique linguistique, assistait au colloque de Strasbourg, à côté de son ami Paul Imbs. On lui prêterait à tort, en cette occasion, un rôle de promoteur de ce qui allait devenir la lexicométrie. Sa conversion n'intervint qu'un peu plus tard, sur le chemin de Besançon, auprès de Quemada, Evrard et Moreau. Mais en 1957, Pierre Guiraud, qui allait publier l'année suivante « Problèmes et méthodes de la statistique linguistique » se trouvait aussi à Strasbourg parmi les intervenants. Certes la statistique ne fut guère évoquée dans les débats (on parlait plutôt de relevés et de dénombrements) et Guiraud choisit de parler d'un sujet un peu moins provocateur : l'argot. Mais les participants au colloque avaient pour la plupart une idée plus confuse encore de l'informatique (le mot n'existait pas alors), que certains confondaient alors – ce n'était pas le cas de Quemada ou de Wagner – avec la mécanographie.

Mais dans les années qui ont suivi le Colloque de Strasbourg, la statistique, présumée un peu plus accessible que l'informatique aux esprits littéraires, bénéficia soudain d'un essor remarquable auquel contribuèrent le Centre d'étude du vocabulaire français de Besançon, l'entreprise du « français fondamental », les travaux d'Evrard en Belgique et du Père Busa en Italie, les publications de Guiraud et, bien sûr, la thèse et le manuel de Muller.

Cela est si vrai que la première publication du TLF a été de nature statistique. C'est en 1971, l'année même où le premier tome du Trésor allait être livré aux lecteurs, que paraissent, avec la signature de Robert Martin, les quatre volumes du *Dictionnaire des fréquences*.

Cinquante ans après on doit reconnaître que la veine statistique dans la mine de Nancy n'a pas été exploitée autant qu'on pouvait l'espérer. Le filon qui devait conduire au trésor n'a pas suscité les vocations attendues, malgré les efforts des directeurs du TLF. Inversement, l'informatique, un peu timide les premières années, a pris un essor spectaculaire dans l'entreprise nancéenne, quand un informaticien de grand talent, Jacques Dendien, a rejoint l'INaLF.

Il n'y a pas concurrence entre l'informatique et la statistique. Les deux sont associées dans les moteurs de recherche, le data mining et la plupart des industries de la langue. Et si l'ATILF rend des services inégalés en matière documentaire, les informations statistiques qu'il distribue sont d'un grand intérêt et donnent à tout le moins les données de base pour les études quantitatives. On se propose d'en donner quelques exemples tirés de la BHVF, du TLFi et de *Frantext*, et peut-être aussi quelques leçons.

Mais avant de procéder à la défense et illustration de la statistique linguistique, il convient, puisque le présent colloque est une rencontre où le témoin va du passé à l'avenir et du bilan au projet, de s'interroger sur les raisons qui expliquent le piétinement d'une discipline scientifique qui s'est pourtant imposée dans la plupart des sciences sociales : sociologie, psychologie, géographie humaine, économie, sciences politiques. Les compteux de mot, pour reprendre une expression qui a cours au Québec, n'ont pas bonne presse, ni chez les linguistes, ni chez les littéraires. Les premiers s'appuient sur des exempliers et trouvent la garantie dans l'attestation d'un fait de langue et non dans sa fréquence. Les seconds se fondent sur les sources, les lectures, la culture et cherchent la garantie dans l'accord avec les jugements d'autrui. En réalité, les réticences s'adressent moins à l'informatique qu'à la statistique. Les vertus domestiques de l'ordinateur ont fini par être reconnues : rares sont les critiques ou écrivains qui écrivent à la main ou qui utilisent encore une machine à écrire. Et depuis qu'*Internet* a gagné en puissance, en extension et en rapidité, les mêmes littéraires sont sensibles aux facilités documentaires que permet le réseau mondial. Ce n'est plus seulement la référence d'un livre qu'on trouve dans le catalogue monstrueux de Google qui tend à se donner l'image de la bibliothèque universelle imaginée par Borges. C'est aussi à l'information primaire, au contenu d'un article ou d'un document, qu'*Internet* donne un accès immédiat. Plus besoin d'attendre, plus besoin de commander, ni même de payer. Certes les moteurs de recherche ne sont pas des critiques avertis : ils mêlent l'ivraie au bon grain. Mais l'indice de notoriété qu'ils utilisent est assez souvent suffisant pour faire apparaître le document pertinent parmi d'autres qui le sont moins. Les détracteurs littéraires de l'informatique et d'*Internet* sont ainsi devenus moins virulents. Ils acceptent les services auxiliaires et ancillaires de la machine, qui réunit les matériaux sur la table de travail. Ils se réservent la part noble de la sélection, de l'exploitation et de l'interprétation. C'est là ce qui compte.

Et pour cela, croient-ils, nul besoin de ceux qui comptent. On ne compte, disent-ils, que ce qui est quantifiable : les prix et les produits, les carottes et les avions. On ne compte pas les idées... Pourtant la démocratie s'exerce en comptant les votes, les opinions, les hommes. Et le jugement littéraire lui-même n'est pas tout à fait dépourvu de compteurs inconscients : beaucoup des jugements qu'on croit qualitatifs sont inspirés par une statistique implicite qui n'avoue pas son nom et qui autorise l'emploi des mots typique, spécifique, caractéristique, si fréquents sous la plume de la critique lorsqu'elle analyse un auteur, un genre ou une époque. Le mot fréquent lui-même relève de cette approche, comme aussi rare, original, banal, courant, cliché, surprise, rupture. Les littéraires parlent d'horizon d'attente quand les statisticiens parlent d'espérance mathématique. L'espérance des uns et l'attente des autres, ce n'est qu'une prévision fondée sur les observations répétées que la conscience enregistre.

Cette réticence des milieux littéraires est d'autant plus regrettable que les données textuelles ont des propriétés très avantageuses que Guiraud avait soulignées, dès 1959, deux ans après le colloque de Strasbourg : « La linguistique est la science statistique type ; les statisticiens le savent bien ; les linguistes l'ignorent encore.¹ » Ces données sont d'abord bon marché ; il suffit d'un traitement de texte pour les enregistrer, ou d'un scanner, ou tout simplement d'une liaison à *Internet* où les textes foisonnent. Elles sont abondantes ; or la statistique se plaît dans les grands nombres et ses conclusions sont moins sûres quand les observations sont en nombre limité parce qu'elles sont chères, ce qui est le cas des enquêtes et des sondages, ou que la nature des choses le veut ainsi, comme il arrive en médecine. Les données textuelles sont faciles à contrôler, à reproduire, à transmettre, à corriger, et se prêtent volontiers à la répétition et aux variations des expériences. Elles sont exemptes de filtrage sélectif et subjectif, et la conscience du chercheur intervient peu à l'entrée. Enfin et surtout les textes sont immédiatement interprétables, sans le truchement des machines. La lecture et la connaissance du texte sont des garants contre les calculs aberrants et les manœuvres avortées. La conscience a au moins une idée de l'ordre de grandeur des résultats attendus et peut les rejeter s'ils sont exagérément erronés ou délibérément triviaux, alors que d'autres disciplines travaillent

¹ Pierre Guiraud, *Problèmes et méthodes de la statistique linguistique*, D. Reidel Publishing Company, Dordrecht-holland, 1959, p. 15.

sans vision directe et se trouvent liées aux instruments, sans possibilité de récuser leur témoignage.

Ce qui est un avantage est aussi un inconvénient. Puisqu'en matière littéraire ou linguistique, on peut se débrouiller sans appareillage, pourquoi s'encombrer de méthodes indirectes et grossières qui obscurcissent l'évidence ou brutalisent les nuances. Pour beaucoup d'esprits, le refus d'accueillir les méthodes quantitatives est une manière de réserver un espace de liberté pour la conscience individuelle. Qu'il y ait au moins un domaine préservé, une réserve naturelle comme il y des parcs du même nom pour la sauvegarde des paysages, un refuge contre l'invasion technologique, où puissent vivre et subsister les disciplines menacées : l'art, le langage, la philosophie, la religion, la musique, la cuisine...

Là s'arrête notre plaidoirie à charge et à décharge. La défense et illustration de la statistique ne passe pas par des discours mais par des exemples. Nous nous proposons d'en donner quelques-uns.

Avec près de 4 000 textes de la littérature nationale, engrangés méthodiquement depuis trente ans, *Frantext* n'a guère d'équivalent dans les autres langues, ni pour l'étendue, ni pour l'homogénéité des données, ni même pour leur accessibilité. La recherche de contextes est de loin la fonction la plus utile et la plus sollicitée de *Frantext*. Certains utilisateurs souhaiteraient certes que la transmission des textes s'ajoute à celle des contextes. Mais ce serait violer les prescriptions du copyright qui n'autorisent qu'un maximum de 300 caractères s'il s'agit d'un texte sous ayant droit. En revanche on n'a pas limité le nombre des contextes restitués pour un mot donné (ou une liste de mots). Les sorties peuvent donc se développer à loisir et cela est favorable à la statistique qui n'aime pas être à l'étroit.

Car l'utilité de *Frantext* ne se réduit pas aux seules opérations documentaires, si sophistiquées soient-elles. Les fonctions statistiques qu'on y trouve ne le cèdent en rien pour la puissance et la portée, et leur exploitation intensive et systématique permet d'atteindre des résultats dont la conscience linguistique, réduite à ses seuls moyens, serait incapable. Ce domaine est toutefois plus technique et moins familier aux populations littéraires qui constituent la clientèle privilégiée de *Frantext*. Et pour ne pas trop effrayer les néophytes, les fonctions statistiques offertes par le logiciel *Stella* ont une simplicité voulue, qui s'arrête aux pourcentages. Cela est suffisant pour donner une idée de la distribution

d'une forme parmi les époques, les écrivains, les textes ou les genres, et ce qui vaut pour une forme peut s'étendre à une constellation lexicale constituée librement autour d'un thème ou d'une construction syntaxique.

Mais trop de simplicité peut conduire à l'erreur d'interprétation et il est dangereux d'accorder une confiance illimitée aux effectifs non pondérés ou aux méthodes trop frustes de pondération que sont les pourcentages et les fréquences relatives. Pour tirer pleinement profit du gisement, il a paru utile d'installer à la sortie de *Frantext* une unité de transformation, qui puisse assurer le traitement quantitatif des matériaux. C'est l'objet du présent logiciel dont le nom THIEF indique assez sa relation de dépendance parasitaire à l'égard de *Frantext* et dont l'écran d'accueil (figure 1) propose deux choix principaux, selon qu'on souhaite être en liaison directe ou différée avec *Frantext*. Dans le premier cas on activera les fonctions rangées verticalement sur la droite de l'écran. Dans le second, on ne s'intéressera qu'aux fonctions disponibles horizontalement au haut de l'écran.

L'exploitation statistique off line

Comme on ne manipule ici que des nombres, en dehors de tout contexte, l'embarras du copyright n'est plus à craindre. Et le traité signé à Genève le 20 décembre 1996 exclut les banques de données des mesures de protection qui entourent la propriété intellectuelle et artistique. Il a donc été possible de livrer au public un dictionnaire de fréquences, qui est issu de *Frantext* par les voies autorisées et qu'on pourrait remettre à jour par les mêmes moyens. L'avantage recherché n'est pas seulement de faire l'économie des liaisons télématiques car, une fois réglé l'abonnement annuel et forfaitaire, le coût pour l'utilisateur est le même, quels que soient le nombre et la longueur des séances de consultation. Ce qu'on recherche surtout, c'est le gain de temps et le confort de l'utilisation que permet tout traitement local. La base locale est accessible aux boutons rangés au haut de l'écran (figure 1). Elle rend compte de l'usage littéraire, les textes dits « techniques » ayant été volontairement écartés. Si elle ne livre pas l'exhaustivité du corpus actuellement disponible, qui s'accroît chaque année, elle n'en représente pas moins l'essentiel des données de *Frantext*, soit 2 376 textes et 117 millions d'occurrences, c'est-à-dire le contenu littéraire de *Frantext* en 1995.



Figure 1. Le logiciel THIEF. Écran d'accueil de la version Windows

	Nb,mots	Nb,Formes	prob,p	prob,q	époque
1	1 719 178	67 014	0,014625	0,985375	1550
2	8 346 862	101 892	0,071006	0,928994	1630
3	6 087 533	69 612	0,051786	0,948214	1692
4	9 380 093	77 841	0,079796	0,920204	1735
5	11 946 384	99 028	0,101627	0,898373	1780
6	11 124 272	98 90	0,094633	0,905367	1820
7	16 184 517	124 845	0,13768	0,86232	1855
8	13 780 168	116 085	0,117227	0,882773	1885
9	8 695 375	98 488	0,073971	0,926029	1910
10	11 361 661	109 218	0,096653	0,903347	1928
11	10 083 262	106 498	0,085777	0,914223	1942
12	8 842 284	112 367	0,07522	0,92478	1960
TOTAL	117 551 589	393 848			

Tableau 2. Limites des 12 tranches

Cette base locale comprend le relevé de toutes les formes du corpus littéraire et des sous-fréquences de chacune dans 12 tranches

chronologiques distinguées du XVI^e siècle à nos jours. Les limites des tranches n'ont pu être établies sur un pied d'égalité, car les textes dépouillés sont très inégalement répartis selon les siècles. Afin d'équilibrer la taille des sous-ensembles, l'empan chronologique a été élargi là où les textes étaient rares, c'est-à-dire au XVI^e siècle, et resserré là où ils abondaient, aux XIX^e et XX^e siècles. La première tranche s'étend ainsi sur un siècle (on l'a représentée par son année médiane : 1550) tandis que les plus proches ne recouvrent guère que deux décennies. Voir tableau 2.

Même ainsi, l'égalité dans l'étendue des tranches n'est pas respectée et les calculs de pondération sont inévitables. Ils s'appuient tous sur les probabilités indiquées dans le tableau précédent. On renvoie le lecteur aux ouvrages de Charles Muller pour tout ce qui concerne les opérations techniques de la statistique linguistique². On s'en fera toutefois une idée suffisante si l'on sait que toute observation réelle (pour un mot donné dans une tranche donnée) est comparée à une fréquence théorique, obtenue par une règle de trois, sur la base de l'étendue respective des tranches. Le résultat de cette comparaison est un nombre négatif ou positif dont le signe indique s'il s'agit d'excédent ou de déficit et dont la valeur absolue mesure l'importance de l'écart (quand l'écart est faible, entre -2 et + 2, le hasard peut être invoqué et l'on doit surseoir à toute conclusion).

Ajoutons qu'une autre base a été extraite des données de *Frantext*, pour rendre compte non plus de la chronologie, mais de l'écriture propre aux écrivains. On a choisi ceux qui étaient les mieux représentés dans *Frantext*, soit 70 écrivains de Rabelais à Gracq, et pour chacun d'eux une fonction puissante de *Stella* a permis de recueillir le dictionnaire des fréquences observées dans son œuvre. Nous ne montrerons ni l'une, ni l'autre de ces deux bases, qui ont été largement exploitées et qui exploitent les fonctions statistiques de notre logiciel HYPERBASE, mais non les fonctions documentaires puisque seules les données numériques y sont accessibles.

Il semble opportun de montrer plutôt comment on pourrait les reconstruire, à partir du corpus enrichi de *Frantext*, ou, mieux encore,

² Charles Muller, *Initiation aux méthodes de la statistique linguistique*, Hachette Université, 1973, et *Principes et méthodes de statistique lexicale*, Hachette Université, 1977. Ces deux ouvrages ont été réédités chez Champion, dans la collection Unichamp.

comment on peut exploiter *Frantext*, en direct, sans s'enfermer dans un corpus figé.

L'exploitation statistique de *Frantext* on line

Les fonctions statistiques apparaissent clairement dans les propositions de *Frantext*, à côté des fonctions purement documentaires. On les a mises en relief dans la figure 3 :

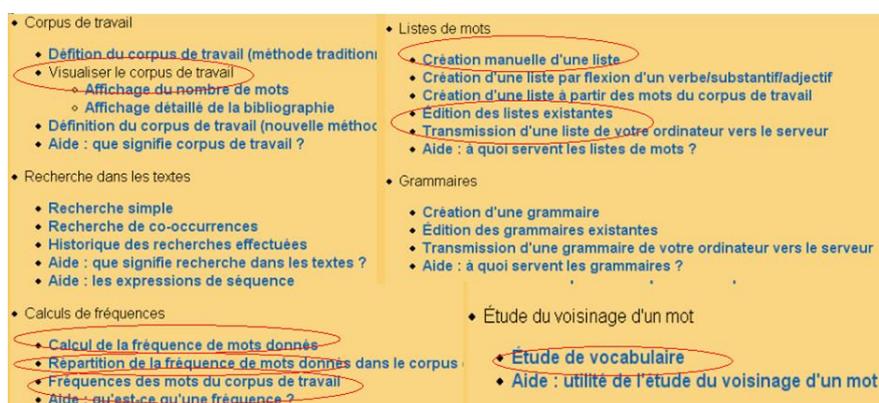


Figure 3 : Les fonctions de *Frantext*

1. L'évolution d'un mot ou d'un groupe de mots

La fonction la plus simple est celle qui rend compte de l'évolution d'un mot (ou d'une liste de mots). Elle suppose comme toute interrogation adressée à *Frantext* que l'on ait d'abord choisi un corpus de travail. On le choisira généralement large, pour permettre au temps de faire son effet. Le bouton à solliciter est le deuxième du groupe « Calculs des fréquences » et s'intitule « Répartition de la fréquence... ». Bien sûr il incite à sélectionner un mot ou une liste mais il permet aussi deux options : l'une s'engage du côté des auteurs, l'autre du côté de la chronologie. C'est ce dernier choix qu'il faut faire en précisant en outre la périodicité que l'on veut, par exemple 20 ou 30 ans.

Si le mot est rare, ou les tranches courtes, il peut se faire que l'effectif étant nul dans une tranche particulière, il n'en soit pas fait mention dans les résultats. Or cette absence peut constituer un déficit intéressant. Le calcul prend donc pour base uniquement les périodes qui contiennent au moins une fois le mot considéré. Les tranches chronologiques qui ont complètement ignoré le mot en question sont

tenues à l'écart du vote et considérées comme abstentionnistes. Ce cas ne se produit guère, si le mot est fréquent ou s'il s'agit d'une liste de mots. Et cela peut se justifier si le mot est de création récente et qu'un déficit dans les siècles anciens est une évidence trop triviale pour être soulignée.

Quand le fichier des résultats a été sauvegardé³, on peut abandonner *Frantext* et le navigateur et solliciter le bouton *Chrono* du menu principal de THIEF. Le résultat final est un histogramme qui apparaît d'abord sur l'écran puis est transmis à l'imprimante, le cas échéant (figure 4).

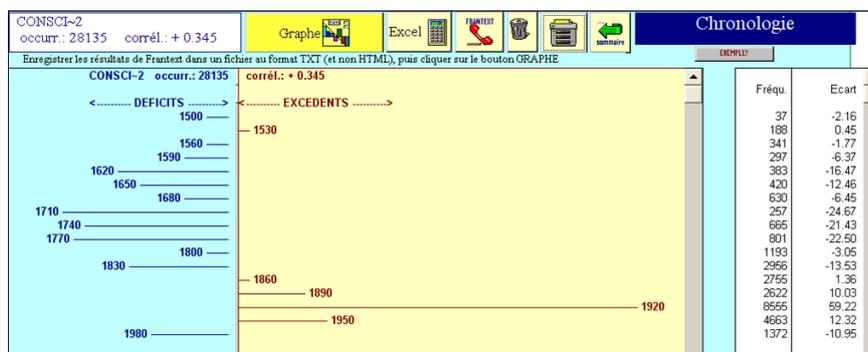


Figure 4. Évolution du mot CONSCIENCE dans le corpus entier de *Frantext* 2008

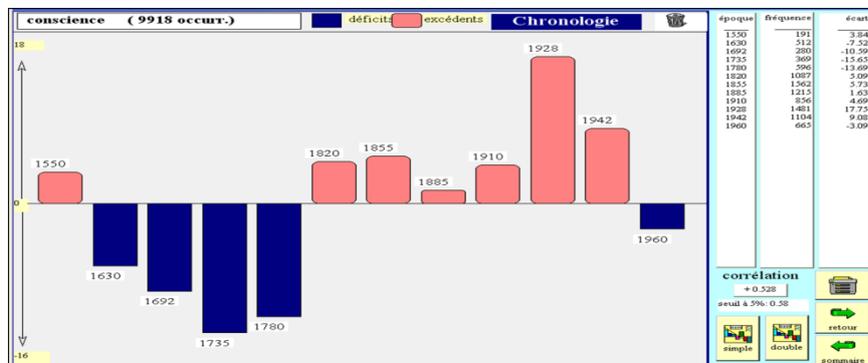


Figure 5. Évolution du mot CONSCIENCE dans le corpus littéraire de *Frantext* 1995

³ Prendre garde à enregistrer les résultats en mode texte. Malheureusement certains navigateurs ignorent dans ce format les tabulations du fichier original, si bien que les nombres étant collés les uns aux autres la lecture en devient impossible. Le navigateur *Firefox* ne présente pas cet inconvénient, non plus que *Camino* dans le monde Mac.

Si l'on préfère une autre présentation, solliciter le bouton « Excel », en précisant que les données sont dans le fichier *chrono.xls*. On peut aussi solliciter la base locale (off line) qui donne un résultat équivalent où la faveur du mot CONSCIENCE doit certes beaucoup à Freud, mais aussi au mouvement spiritualiste qui l'a précédé et dont Bergson est la figure de proue (figure 5). Mais les choix ne sont plus modifiables, ni dans la composition du corpus ni dans la périodicité.

2. La répartition d'un mot ou d'un groupe de mots chez les écrivains

Il s'agit là de la deuxième option du programme « Répartition » de *Frantext*. La procédure d'interrogation et de sauvegarde est la même. Mais le résultat est sensiblement plus long, ayant autant de lignes que d'écrivains recensés pour le mot en question, même si ce mot ne figure qu'une fois dans son œuvre. On en donne ci-dessous un extrait, relatif au mot AME (figure 6).

Fréquence absolue totale : *59303*

Fréquence relative maximale : *4977* chez *RODENBACH Georges*

Diagramme des fréquences relatives

Échelle : un astérisque représente une fréquence relative de *200 millionième(s)

		freq. abs.	freq. rel.	
1	RODENBACH Georges *****	108	4977	
2	SAMAIN Albert *****	84	3807	
3	LESPINASSE Julie de	556	3374	*****
4	DESBORDES-VALMORE Marceline	304	3261	*****
5	NOAILLES Anna de	43	2953	*****
6	POULET Georges	19	2533	*****
7	GUÉRIN Charles	53	2458	*****
8	GUÉRIN Maurice de	509	2407	*****
9	GUÉRIN Eugénie de	586	2119	*****
10	MONOD Henri	70	1959	*****
11	BRIZEUX Auguste	48	1927	*****
12	DIERX Léon	80	1817	*****
13	MAETERLINCK Maurice	259	1762	*****
14	MAINE DE BIRAN	608	1715	*****
15	JOUBERT Joseph	173	1639	*****
16	LAFORGUE Jules	40	1576	*****
17	MORÉAS Jean	83	1564	*****
18	SULLY PRUDHOMME Armand	85	1536	*****

Figure 6. Le résultat proposé par *Frantext* pour le mot AME

Le résultat final est un histogramme qui apparaît d'abord sur l'écran puis est transmis à l'imprimante, si on le désire. Noter que la place est

trop exiguë pour représenter la série dans son intégralité. Dans certains cas, c'est plus d'un millier d'auteurs qui sollicitent chacun un « bâton » de l'histogramme. Le programme permet de ne retenir que les plus importants, c'est à dire les mieux représentés dans le corpus. Le seuil proposé pour l'étendue est de 200 000 occurrences, mais l'utilisateur peut déplacer ce seuil vers le haut ou le bas.

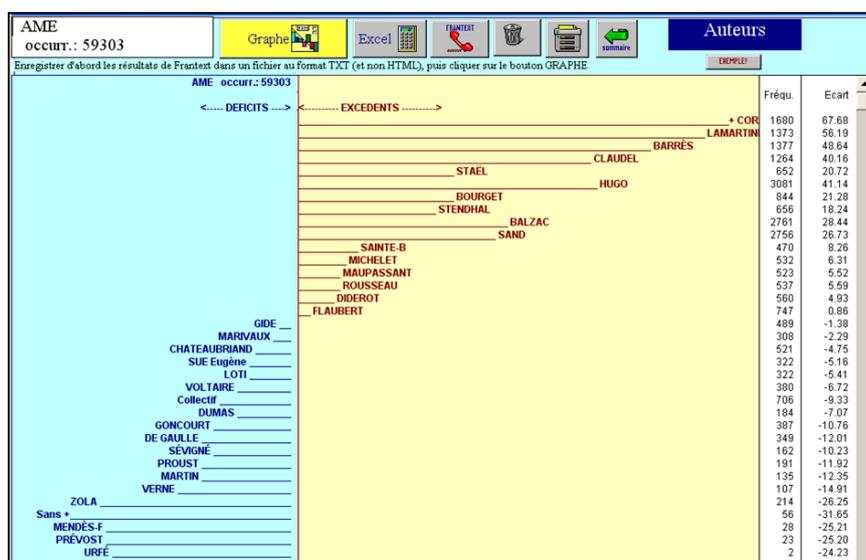


Figure 7. Le graphique du mot AME dans *THIEF*

Le classement est fondé sur les fréquences relatives. La décroissance est régulière. Le même classement est repris dans l'histogramme proposé par THIEF (graphique 7). Mais il est contesté par l'écart réduit, qui suit la tendance générale mais s'en écarte dans les détails. Si on teste le supplément d'âme, l'écart réduit donnerait l'avantage à Hugo sur Mme de Staël et à Balzac sur Stendhal. Les fréquences relatives ont l'inconvénient d'ignorer la taille des observations et d'établir l'égalité entre 20/100 et 20 000/100 000, ce qui est fort différent en terme de probabilité. Là-dessus on renvoie le lecteur à l'avertissement de Charles Muller qui dès 1968 opposait l'impasse et le bon chemin⁴. Comment retrouver le bon

⁴ « Cette présentation par pourcentages a une petite allure statistique, mais il faut résister à la tentation de l'appliquer, car c'est une impasse qui n'aboutit à aucune conclusion. La bonne façon n'est d'ailleurs pas plus compliquée [...] Une simple règle de trois conduit à une valeur théorique, à partir de laquelle on obtiendra un écart absolu avec la valeur observée ou réelle. », *Initiation à la statistique linguistique*, 1968, p 44.

chemin ? Prenons la ligne relative à Albert Béguin qui indique une fréquence absolue de 330 et une fréquence relative de 1 488. Sachant que la fréquence relative a été obtenue en divisant la fréquence du mot par la taille du texte (résultat multiplié par 1 000 000), on retrouve la taille du texte en faisant l'inverse, soit : $(330 / 1\,488) \times 1\,000\,000 = 221\,774$. On procède ainsi pour toutes les lignes, en faisant le total pour la fréquence du mot et la taille du corpus entier, soit : fréquence de âme : 59 303, taille du corpus : 202 633 773. On peut alors calculer la fréquence théorique par la règle de trois : $59\,303 \times (221\,774 / 202\,633\,773) = 64,90$

puis l'écart absolu : $330 - 64,90 = 265,10$

et enfin l'écart réduit (formule simple) : $z = 265,10 / \sqrt{64,90} = 32,90$

3. Constitution et traitement des tableaux

Cette fonction est plus riche et plus puissante que les précédentes. Mais elle requiert une consultation plus longue de *Frantext*. Car il faut de façon répétitive interroger *Frantext* en lui proposant la même liste de mots mais en variant la composition du corpus. Le choix de la liste peut être fait en conversationnel et se maintient tant que dure la communication. Mais on peut l'établir dans un fichier local qu'on transmet à *Frantext* au moment de la consultation. Quand on a épuisé la liste des auteurs (ou des textes, des genres ou des époques), on dispose d'autant de fichiers, qu'on va réunir dans un tableau (tableau 8), chaque fichier servant à remplir une colonne.

	ALEM	BALZ	CLAU	DIDE	FLAU	GIDE	GRAC	HUGO	MARI	PROU	RIMB	SART	STEN	VERN	VIGN	VOLT	ZOLA
ils	524	9565	2027	3517	5261	2920	969	8080	1725	2634	66	2295	1337	2939	928	555410178	60519
je	495	46084	18600	17683	32501	34026	5406	42216	21346	15127	424	8155	9929	2484	4902	1047014094	283942
leur	682	8806	1383	2174	3037	2153	1101	5988	1071	2552	66	776	1120	1968	905	3825 5207,	42814
lui	1255	22796	5535	5464	7208	6862	2482	13995	5075	7095	55	1747	5876	2792	1650	683018962,	115679
m'	109	10405	3745	4032	8720	9060	1445	8280	5361	4347	72	1593	2063	369	1135	2440 2954,	68110
me	134	15891	5559	6084	10964	12458	2439	11024	9535	7312	126	2428	4160	566	1978	3049 3938,	95545
moi	123	9352	6347	3476	6148	5048	1395	9824	4017	3734	61	1557	2077	508	1190	1951 3587,	60395
nous	1661	17876	10025	5299	8766	9493	2404	16140	4814	6676	188	2021	2092	2762	1341	4991 5901,	102450
on	1104	13579	4812	9129	10679	8525	4304	26133	5064	6820	138	2145	4154	4117	1972	1176014579,	129014
se	1184	33073	5691	6941	11991	8963	4894	21127	2917	6051	142	3268	6443	7271	2213	703927885,	157093
soi	8	302	184	179	288	406	132	361	41	232	0	50	124	23	61	121 165,	2677
t'	47	2226	1307	343	3503	1209	271	3295	487	147	15	600	187	53	208	474 1265,	15625
te	72	3615	1946	385	3363	1410	353	3620	772	251	34	797	390	94	217	500 1842,	19651
toi	52	1963	2377	256	2655	1101	262	3550	351	134	23	547	203	69	184	559 1319,	15605
tu	119	8918	6002	1365	9080	3657	1284	9311	1750	730	87	2955	726	322	601	1398 6365,	54670
vous	488	49764	12438	16415	16322	8443	1402	37687	19743	6183	89	3214	5346	3034	2774	1074112975,	207058
	8057	254215	87978	82742	140486	115734	30543	220601	83069	70025	1586	34148	46217	29371	22157	71702131216,	1429847

Tableau 8. Un tableau de contingence

À l'intersection d'une ligne et d'une colonne, on lit alors la fréquence du mot x dans le corpus y . Dès lors toutes les manipulations qu'on peut faire avec un tableau de contingence sont possibles. Il suffit de

s'appuyer sur le total général et les totaux marginaux de ligne et de colonne pour produire fréquences théoriques et écarts réduits, et conséquemment histogrammes, analyses arborées et analyses factorielles (figure 9). La page Tableau de THIEF offre tous ces outils traditionnels du traitement statistique.

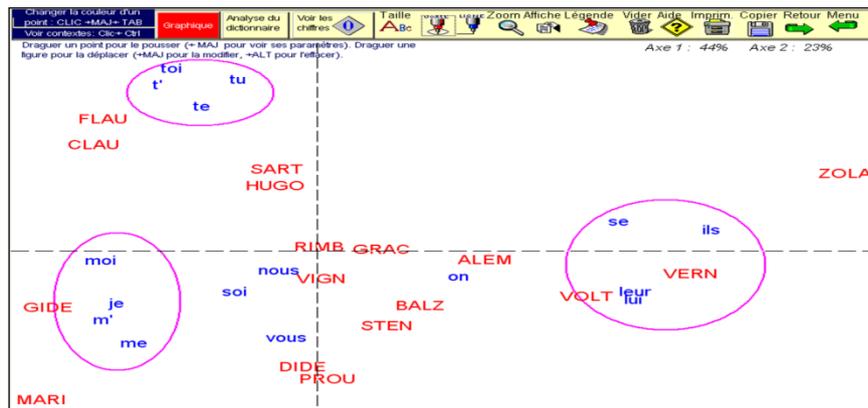


Figure 9. Analyse factorielle du tableau 8

4. Le vocabulaire spécifique d'un texte

Cette proposition de *Frantext* est puissante et précieuse, puisqu'elle permet en quelques secondes de connaître et d'enregistrer la liste alphabétique de tous les mots d'un corpus, si gros soit-il, chaque forme étant suivie de sa fréquence. On appelle cela un dictionnaire de fréquences. Les fréquences sont absolues et livrées brutes. Mais il est facile d'en tirer parti, si on dispose d'un élément de comparaison. THIEF propose un tel modèle, qui est lui-même issu de *Frantext* et qui regroupe les textes littéraires, soit 117 millions de mots (cela correspond à la base locale évoquée plus haut). Prenons pour exemple le Roman de Gracq, *Un Balcon en forêt*, actuellement au programme de l'agrégation. Si un agrégatif se trouve dans la salle, il sera peut-être intéressé par la figure 10, qui dresse la liste des spécificités de ce texte, en ordre hiérarchique et alphabétique. Il n'y a pas lieu d'admirer la position avantageuse des noms propres (Grange, Meuse, Mona) liée à des privilèges évidents. Mais que le premier terme significatif soit le mot «forêt» n'est pas dû au hasard. À lire la liste de ces descripteurs quantitatifs, on devine le sujet et l'atmosphère, comme s'il s'agissait d'une fiche signalétique.

Enregistrer d'abord les résultats de Frantext dans un fichier au format TXT (et non HTML), puis cliquer sur le bouton CALCUL

Vocabulaire spécifique d'un texte		Excédents				Déficits			
		excédents	écart	texte	corpus	déficits	écart	texte	corpus
EXEMPLE!		grange	195.08	309	1470	je	-22.58	84	396110
		mieuse	84.95	73	431	est	-17.66	153	342841
		mona	62.52	54	251	j'	-14.79	24	157060
		forêt	60.79	128	2471	.	-14.33	2025	1652388
		fortin	58.38	18	56	nous	-13.61	24	136975
		béton	48.41	30	223	a	-13.35	57	166682
		taillis	45.67	40	440	...	-12.42	123	212559
		dégel	43.09	15	71	que	-11.94	404	431159
		camionnette	35.36	18	150	vous	-11.91	86	171164
		clairière	32.16	21	244	ai	-11.22	16	92692
		sentait	31.22	105	5582	ne	-11.04	216	267396
		glissait	26.45	29	662	pas	-10.22	283	305816
		capitaine	25.57	56	2451	tu	-9.92	29	69602
		branches	24.85	48	1931	mon	-8.84	32	79885
		cavalerie	23.66	25	611	elle	-8.13	235	235834
		pensa	23.65	45	1866	bien	-8.07	61	97700
		belgique	23.28	24	582	ce	-7.78	233	229220
		frontière	21.21	26	804	pour	-7.64	199	201936
		toit	21.00	36	1511	suis	-6.28	17	41150
		la	21.00	2119	801535	fait	-5.95	45	63189
		julia	20.91	18	408	ou	-5.73	49	64607
		jusqu'	20.04	72	5755	sont	-5.71	25	44055
		neige	19.48	44	2488	faire	-5.40	36	51125
		songea	18.32	21	697	qu'	-5.37	385	300597
		chemin	18.04	71	6573	n'	-5.18	249	205106
		lourde	17.74	35	1911	dire	-5.11	22	36997
		bois	17.60	86	9445	si	-4.95	109	103801
		roule	17.43	77	7935	été	-4.75	19	31993
		songeait	17.30	27	1236	peut	-4.65	35	44793
		guerre	17.04	93	11254	non	-4.55	32	41658
		carré	16.14	20	796	être	-4.52	74	73938
		fit	15.75	103	14703	avoir	-4.35	19	29466
		montait	15.65	30	1775	ça	-4.33	55	58150
		lieutenant	15.31	23	1130	ont	-4.26	19	28930
		dernière	15.26	87	11631	mais	-4.18	197	157506
		était	14.79	484	148707	voir	-4.09	16	25406
			14.46	2003	130220	tous	-3.92	28	43524

Figure 10. Le vocabulaire spécifique d'*Un balcon en forêt*

5. Étude du voisinage d'un mot

Cette fonction de *Frantext* ressemble à la précédente et livre une liste semblable dotée de fréquences brutes. Mais les mots retenus ne sont pas ceux d'un texte suivi, mais d'un agrégat de contextes, réuni autour d'un mot-pôle (ou d'une liste de mots). Il appartient au chercheur ici comme partout de définir le corpus, puis de proposer un mot (ou une liste de mots) que le logiciel *Stella* doit rechercher dans le corpus. *Stella* ne montre pas le monceau de contextes qu'il a trouvés mais en restitue le contenu sous la forme d'un dictionnaire de fréquences. On s'attend que ces fréquences soient liées par quelque raison, thématique ou syntaxique, voire métrique, au mot-pôle, mais la fréquence d'un mot dépend aussi de lui-même⁵, quel que soit le contexte. Quel que soit le mot-pôle, il est probable que les mots HOMME, JOUR et FEMME⁶, que l'on voit presque à

⁵ Pierre Guiraud avait imaginé que la fréquence était une propriété attachée à chaque mot de la langue. On lui a objecté que cette propriété était variable et dépendait de la situation de discours. Il n'en reste pas moins que la disponibilité n'est pas la même pour tous et qu'il est rare qu'un mot soit premier dans un discours et dernier dans un autre. On voit souvent les mêmes au premier rang.

⁶ Ce sont les trois substantifs que *Frantext* place en tête dans le corpus littéraire.

6. Recherche de contextes

C'est là la voie royale de *Frantext*, nullement cataloguée comme statistique. Elle permet pourtant une démarche quantitative lorsqu'on assemble dans un fichier les contextes qu'elle restitue. Un tel fichier a nécessairement été créé dans la fonction précédente, qui explore le voisinage d'un mot. On fait donc ici un pas en arrière en demandant à la machine locale de refaire ce qu'avait fait le serveur. Il en ressort une liste de spécificités comme précédemment. Mais cette fois, comme on dispose du texte, les investigations peuvent aller plus loin. Au lieu de calculer seulement la liaison de chaque mot avec le pôle unique, on peut étudier les relations que chacun des mots entretient avec tous les autres. On rejoint là la démarche très générale où l'étude systématique des cooccurrences aide à trouver des solutions pour des problèmes divers : le data mining, les candidats-termes, la traduction assistée, le résumé automatique, le traitement des profils...

Lecture des contextes reçus de Frantext		Thème	CO
Résultat 1 (Texte sous droits) L669/ GRACQ Julien /Au château d'Argol/1938 Pages 10-16 / AVIS AU LECTEUR qu' elles ont toujours inépuisablement versé sur lui. 1938. quoique la campagne fût chaude encore de tout le soleil de l' après-midi, *Albert s' engagea sur la longue route qui conduisait à *Argol. Il s' abrita à l' ombre déjà grandie des aubépines et se mit en chemin. Il voulait se donner une heure encore pour savourer l' angoisse du hasard. Il	Extrait 1	Voir la suite?	
Résultat 2 (Texte sous droits) L669/ GRACQ Julien /Au château d'Argol/1938 Pages 19-21 / ARGOL y entendait couler des ruisseaux invisibles, mais *Albert fut frappé par la rareté et la triste monotonie du chant des oiseaux. Une hauteur toute proche et parallèle à la route arrêtait la vue de ce côté : là des pins parasols allongés en une mince ligne sur la crête contre le soleil couchant semblaient souligner de leur ramure horizontale et élégante			
Résultat 3 (Texte sous droits) L669/ GRACQ Julien /Au château d'Argol/1938 Pages 20-22 / ARGOL dressait à l' extrémité de l' éperon rocheux que venait de côtoyer *Albert. Un sentier tortueux y conduisait- impraticable à toute voiture -et s' embranchait à gauche de la route. Il serpentait quelque temps dans une étroite prairie marécageuse, à travers laquelle *Albert entendit le plongeon précipité des grenouilles sur son passage. Puis le			
Résultat 4 (Texte sous droits) L669/ GRACQ Julien /Au château d'Argol/1938 Pages 21-23 / ARGOL des dangers ! Les merlons de cette puissante tour ronde, faite de dalles épaisses de granit, se profilaient toujours juste au-dessus de la tête du voyageur engagé dans sa route pénible, et rendaient plus frappante la vitesse des lourds nuages gris qui les débordaient à chaque seconde avec une rapidité sans cesse accrue. à l' instant où *Albert			

Tableau 13. Les contextes du mot ROUTE dans l'œuvre de Gracq (extrait)

Quand le traitement des cooccurrences est achevé, plusieurs outils permettent d'appréhender le tableau triangulaire où chaque mot est relié à chacun des autres par un indice de proximité, à l'image des atlas où la distance de ville à ville se lit dans un triangle analogue. Les ressources habituelles s'appliquent ici : histogrammes, analyses arborée et

factorielle. Nous y ajouterons un graphe où la ROUTE est un carrefour où se rejoignent les compagnons de ROUTE ; les arcs sont en rouge quand il s'agit d'une liaison avec le pôle, en bleu lorsqu'il s'agit d'une ligne intérieure ou transversale. La force du trait (gras, maigre ou pointillé) indique la force de la liaison. Les liaisons de la ROUTE⁷ sont bien celles qu'on attendait chez un écrivain-arpenneur qui a beaucoup marché et circulé dans la campagne française et qui est l'un des premiers paysagistes de notre littérature.

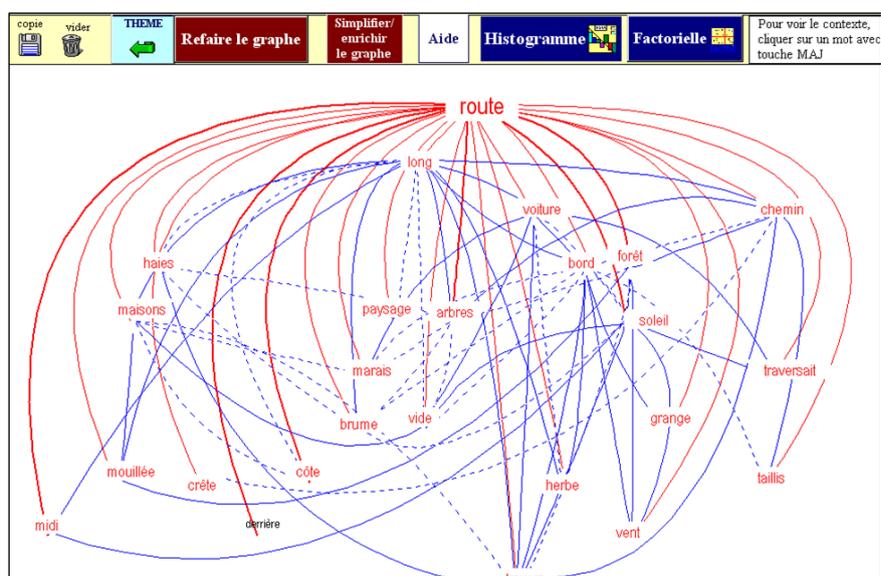


Figure 14. Le graphe de la ROUTE

Conclusion

1- Le graphe qu'on vient de montrer rend possible le retour au texte. Il suffit de cliquer sur un mot pour voir apparaître les contextes où on l'emploie. Un tel retour est malheureusement impossible pour la plupart des fonctions statistiques disponibles dans *Frantext*, qui ne transmettent que des nombres. Si *Frantext* délivrait aussi les textes, les ateliers spécialisés de statistique auraient un matériau plus fruste mais plus riche. Or là où le copyright n'est pas un droit opposable, cette libre transmission

⁷ Un de ce que Gracq qualifie de « fragment » porte ce nom. Et son dernier livre a pour titre « Carnets du grand chemin ».

du trésor serait très souhaitable. Au reste il y a des précédents, dont ont bénéficié l'ARTFL de Chicago, la Bibliothèque Nationale et quelques privilégiés, dont je suis.

2- Dans son état actuel les fonctions statistiques de *Frantext* ne s'appliquent qu'aux formes graphiques. On souhaiterait qu'elles s'étendent aux lemmes, puisqu'aussi bien les textes ont été, dans leur majorité, étiquetés et lemmatisés.

3- On aimerait qu'un véritable atelier statistique soit ouvert à Nancy et que le serveur puisse offrir des traitements statistiques évolués, et pas seulement des fréquences et des pourcentages. Un tel atelier irait au bout du traitement, y compris jusqu'aux graphes comme ceux que délivre le Crisco de Caen ou le CNRTL de Toulouse. La situation actuelle, où les produits échangés sont demi-finis, engendre trop de contraintes et de corrections, dès que change le format des données.

4- On réclame enfin que dans le domaine des données chiffrées le copyright des éditeurs soit abandonné et que l'abonnement ne soit plus nécessaire, puisqu'à aucun moment le texte n'est communiqué.

5- Ce sont là des projets à court terme dans le développement des bases constituées à Nancy. Ils n'excluent pas des projets plus ambitieux que je laisse aux plus jeunes le soin d'évoquer. Je me suis laissé dire qu'un projet se préparait qui porterait sur l'étiquetage non plus seulement grammatical, mais sémantique des textes. Voilà une recherche d'avenir, à côté de laquelle ce que je viens de raconter n'est qu'un jeu d'enfant.