

Ludovic Lebart

Directeur de Recherche au CNRS,
Ecole Nationale Supérieure
des Télécommunications

André Salem

Ingénieur à l'Ecole Normale Supérieure
de Fontenay-Saint-Cloud

Statistique textuelle

Préface de Christian Baudelot

Professeur à l'Ecole Normale Supérieure

Ouvrage publié initialement par Dunod en 1994

Préface

Et le Verbe s'est fait Nombre...

Il y a dans l'activité qui consiste à traiter les mots comme des nombres - opération de base de la statistique textuelle - un a priori qui ne manquera pas d'apparaître à certains comme outrageusement réducteur voire même sacrilège. Surtout si l'on en croit Victor Hugo : *Car le mot, c'est le Verbe, et le Verbe c'est Dieu...*

Il suffit de lire ce livre et surtout d'en appliquer les principes à ses propres enquêtes pour se convaincre du contraire. Avec ses graphes d'analyse factorielle, J.P. Benzécri a rendu les individus à la statistique : longtemps ignorés à force d'être confondus dans de vastes agrégats ou pulvérisés dans des formules inférentielles qui s'intéressent d'abord aux relations entre des grandeurs abstraites (revenu et consommation, salaire et diplôme...), les individus effectuent leur rentrée sur la scène statistique sous la forme de points dans un nuage. Les positions respectives qu'ils occupent au sein de ce nuage démontrent d'abord qu'ils diffèrent tous les uns des autres. Les distances et les proximités qu'ils entretiennent avec les modalités des variables considérées permettent ensuite de comprendre en quoi chacun diffère de l'autre : par ses goûts, ses opinions politiques, son âge, son sexe, la marque de sa voiture, la profession de son père... mais la statistique est encore une histoire sans parole.

L'une des contributions majeure de la statistique textuelle est précisément d'animer tous ces graphes en donnant la parole à chacun de ces individus. Grâce à Lebart et Salem, les fameux points-individus ne sont plus muets, ils parlent. Vole alors en éclats la traditionnelle mais artificielle distinction entre le quantitatif et le qualitatif. Les méthodes ici présentées permettent de mettre en relation les propriétés sociales ou personnelles des individus telles que les saisit l'enquête statistique avec les textes par lesquels ces mêmes individus répondent aux questions qu'on leur pose sans en réduire le moins du monde l'information. Les nuances les plus subtiles de l'expression sont conservées : le singulier et le pluriel, la majuscule et la minuscule, l'usage du "je", du "on", du "nous". La formule le dit bien : *s'exprimer* c'est d'abord se livrer soi-même au-dehors. Chaque forme lexicale tire alors son sens d'un triple registre : celui que lui donne celui qui la prononce, celui que lui confère la place qu'elle occupe dans l'espace dessiné par toutes les autres formes lexicales énoncées par le même individu, celui, enfin, qu'elle tient de la place qu'elle occupe dans l'espace dessiné par toutes les autres

formes énoncées par tous les autres locuteurs. Le sens jaillit des différences de profil.

Cet ouvrage a le mérite de déborder largement le cadre de l'analyse de contenu ou du traitement statistique des questions ouvertes dans les enquêtes. Il fait le point sur l'état de développement d'un chantier particulièrement foisonnant depuis dix ans. Il expose les dernières découvertes. Elles sont nombreuses et riches d'application dans les domaines les plus divers : stylométrie, recherche documentaire, modèles prévisionnels. Comment attribuer un texte à un auteur ou à une période ? Combien d'auteurs ont contribué à la rédaction du livre de la Bible attribué au prophète Isaïe ? Peut-on comparer des comportements exprimés dans des textes écrits dans des langues différentes sans les traduire ni les coder ?

C'est souvent aux confins des disciplines instituées que l'invention scientifique est la plus féconde. Lorsque deux statisticiens tout particulièrement sensibilisés aux problèmes que l'on rencontre dans les sciences humaines se réunissent autour d'un ordinateur pour élaborer les principes et les outils d'une statistique textuelle, ils occupent le cœur d'un carrefour scientifique vers lequel convergent tout naturellement des linguistes, d'autres statisticiens bien sûr mais aussi les spécialistes d'analyse du discours, d'analyse de contenu, d'analyse des textes littéraires, de recherche documentaire et d'intelligence artificielle. A ce noyau dur de producteurs de théories et d'outils est venu petit à petit s'agréger un univers polyglotte d'utilisateurs aux formations diverses : sociologues, littéraires, stylomètres, historiens, géographes, politologues, médecins, éthologues, psychologues, publicitaires, etc.

On peut savoir gré à l'ouverture d'esprit des deux auteurs (et de leurs associés !), à leur générosité intellectuelle et humaine pour avoir su accueillir autour de leur disque dur un nombre croissant de producteurs et d'utilisateurs dont ils ont souvent stimulé l'inventivité. Il suffit pour s'en convaincre de feuilleter les actes des deux journées internationales qu'ils ont suscitées, avec d'autres, à Barcelone en 1990 et à Montpellier en 1993. Ou de goûter, chez soi, le charme inattendu de nouveaux logiciels.

Au-delà de la collection de principes et d'outils statistiques présentés dans les pages qui suivent, n'oublions pas que la nature même de la matière travaillée - le texte - confère à l'entreprise des dimensions à la fois culturelles, internationales et universelles car comme le disait si bien Victor Hugo ...

Christian Baudelot

AVANT-PROPOS

Cet ouvrage s'adresse à ceux qui, pour leurs recherches, leurs travaux d'études, leur enseignement, doivent décrire, comparer, classer, analyser des ensembles de textes. Il peut s'agir de textes littéraires, scientifiques (bibliométrie, scientométrie, recherche documentaire), économiques, sociologiques (réponses aux questions ouvertes dans des enquêtes socio-économiques, entretiens divers en marketing, psychologie appliquée, pédagogie, médecine), de textes historiques, politiques...

On a tenté de faire le point sur les développements de la *statistique textuelle*, domaine de recherche vivant dont les contours exacts sont difficiles à établir tant est large l'éventail des disciplines concernées, et aussi celui des applications possibles. Les chapitres qui suivent voudraient, tout en présentant l'acquis de ce champ disciplinaire, témoigner de cette richesse d'approches, de méthodes et de domaines.

L'ouvrage reprend, en intégrant des développements récents, certains exemples du manuel *Analyse statistique des données textuelles* publié par les mêmes auteurs en 1988. Le champ des applications précédemment limité aux traitements de *questions ouvertes* a été considérablement élargi de même que l'éventail des méthodes proposées. L'ensemble, profondément remanié, inclut de nouveaux chapitres qui traitent des structures a priori et de l'analyse discriminante textuelle, thèmes qui dépassent largement l'optique essentiellement descriptive de l'ouvrage antérieur.

Plusieurs lectures devraient être possibles selon la formation du lecteur, et selon notamment ses connaissances en mathématique et statistique. Une lecture technique, complète, pour une personne ayant dans ces matières une formation équivalente à une maîtrise de sciences économiques, aux écoles d'ingénieurs ou de commerce. Une lecture pratique, d'utilisateur, pour les personnes spécialisées dans les divers domaines d'application potentiels.

Les démonstrations strictement mathématiques ne figurent pas dans le texte. On renvoie à chaque fois le lecteur curieux d'en connaître les détails à des publications ou ouvrages plus spécialisés lorsque ceux-ci sont facilement accessibles. En revanche, la part belle est faite à la définition des concepts, à la mise en oeuvre des procédures, aux règles de lecture et d'interprétation des résultats. Le glossaire en fin d'ouvrage aidera le lecteur à préciser le contenu des notions ou des conventions de notation les plus importantes.

L'ensemble doit beaucoup à des collaborations et des cadres de travail divers : au sein du département Economie et Management, de l'Ecole Nationale Supérieure des Télécommunications (Télécom Paris) et de l'URA820 du Centre National de la Recherche Scientifique (Traitement et Communication de l'Information) de cette même Ecole ; au sein du Laboratoire "Lexicométrie et textes politiques", URL 3 de l'Institut national de la langue française (INaLF) et de l'Ecole Normale Supérieure de Fontenay-Saint-Cloud.

Nous remercions également les autres chercheurs ou professeurs auprès desquels nous avons puisé collaboration et soutien, ou simplement eu d'intéressants débats ou discussions. Citons, sans être exhaustif, C. Baudelot (ENS, Paris), M. Bécue, (UPC., Barcelone), L. Benzoni (Télécom Paris), E. Brunet (INaLF, Nice), S. Bolasco (Univ. de Salerne), L. Haeusler (Cisia, Paris), G. Hébrail (EDF, Clamart), D. Labbé (CERAT, Grenoble), A. Lelu (Univ. Paris VIII), M. Reinert (Univ. Toulouse Le Mirail).

L. L., A. S.
Paris, Janvier 1994

Sommaire

Introduction	7
Chapitre 1 : Domaines et problèmes	11
1.1 Approches du texte	11
1.1.1 Le courant linguistique	12
1.1.2 Analyse de contenu	13
1.1.3 Intelligence artificielle	14
1.2 Les rencontres de la statistique et du texte	15
1.2.1 Les premiers travaux	16
1.2.2 Les banques de données textuelles	17
1.2.3 La recherche documentaire	18
1.3 Approche statistique du texte	18
1.3.1 La chaîne de traitement	19
1.3.2 Connaissances internes et externes	20
1.3.3 Une méta-information exceptionnelle	21
1.4 Des textes particuliers : les questions ouvertes	23
1.4.1 Les questions ouvertes : un outil de recherche	24
1.4.2 Questions ouvertes et questions fermées	25
1.4.3 Quand utiliser les questions ouvertes ?	27
1.4.4 Traitement pratique des réponses libres	28
1.4.5 Les regroupements de réponses	30
Chapitre 2 : Les unités de la statistique textuelle	33
2.1 Le choix des unités de décompte	33
2.1.1 Le texte en machine	35
2.1.2 Les dépouillements en formes graphiques	35
2.1.3 Les dépouillements lemmatisés	36
2.1.4 Les dépouillements à visée "sémantique"	38
2.1.5 Très brève comparaison avec d'autres langues	40
2.2 Segmentation et numérisation d'un texte	42
2.2.1 Numérisation sur le corpus <i>Enfants</i>	44
2.2.2 Le corpus P	45
2.3 L'étude quantitative du vocabulaire	46
2.3.1 Fréquences, gamme des fréquences	46
2.3.2 La loi de Zipf.	47
2.3.3 Mesures de la richesse du vocabulaire	49
2.4 Documents lexicométriques	51
2.4.1 Index d'un corpus	52
2.4.2 Contextes, concordances	53
2.4.3 L'accroissement du vocabulaire	55

2.4.4	Partitions du corpus	56
2.4.5	Tableaux lexicaux	57
2.5	Les segments répétés	58
2.5.1	Phrases, séquences	59
2.5.2	Segments, polyformes	60
2.5.3	Quelques propriétés relatives aux segments	62
2.6	Les inventaires de segments répétés	63
2.6.1	Inventaire alphabétique des segments répétés	64
2.6.2	Inventaire hiérarchique des segments répétés	66
2.6.3	Inventaires distributionnels des segments répétés	68
2.6.4	Tableau des segments répétés	69
2.7	Recherche de cooccurrences, quasi-segments	70
2.7.1	Recherche autour d'une forme-pôle	70
2.7.2	Recherches de cooccurrences multiples	72
2.7.3	Quasi-segments	72
2.8	Incidence d'une lemmatisation sur les comptages	73
2.8.1	Le corpus <i>Discours</i>	73
2.8.2	Principales caractéristiques quantitatives	75
Chapitre 3	: L'analyse des correspondances	79
3.1	Principes de base de l'analyse des données	80
3.2	L'analyse des correspondances	81
3.2.1	Bref historique	81
3.2.2	L'analyse des correspondances exposée à partir d'un exemple simple	82
3.2.3	Validité de la représentation	89
3.2.4	Variables actives et illustratives	92
3.3	Analyse des correspondances multiples	98
3.3.1	Structure de base d'un échantillon d'enquête	100
3.3.2	Validité de la représentation	105
3.3.3	Positionnement des variables illustratives	106
Chapitre 4	: La classification automatique des formes et des textes	111
4.1	Rappel sur la classification hiérarchique	112
4.1.1	Le dendrogramme	113
4.1.2	Coupures du dendrogramme	115
4.1.3	Adjonction d'éléments supplémentaires	116
4.1.4	Filtrage sur les premiers facteurs	116
4.2	Classification des éléments d'un tableau lexical	117
4.2.1	Classification des formes	117
4.2.2	Classification des textes	120
4.2.3	Remarques sur les classifications de formes	123
4.3	La Classification des fichiers d'enquête	127

4.3.1	Les algorithmes de classification mixte	128
4.3.2	Séquence des opérations	130
4.3.3	Exemple d'application	130
Chapitre 5 : Typologies, visualisations		135
5.1	Analyse des correspondances sur tableau lexical	137
5.1.1	Les tableaux lexicaux de base	137
5.1.2	Les tableaux lexicaux agrégés	138
5.1.3	Seuil de fréquence pour les formes	139
5.1.4	Présentation de l'exemple	139
5.1.5	Construction du tableau lexical agrégé	139
5.1.6	Analyse et interprétation du tableau lexical	144
5.2	Les noyaux factuels	148
5.3	Analyse directe des réponses ou documents	152
5.3.1	Comment interpréter les distances ?	152
5.3.2	Analyse du tableau clairsemé T	153
5.3.3	Exemple d'application	154
5.4	Analyse des correspondances à partir d'une juxtaposition de tableaux lexicaux	160
5.5	Analyse des correspondances à partir du tableau des segments répétés	162
Chapitre 6 : Eléments caractéristiques, réponses ou textes modaux		171
6.1	Formes caractéristiques, spécificités	172
6.1.1	Le calcul des spécificités	172
6.1.2	Un exemple de calcul des spécificités	177
6.1.3	Liste des formes spécifiques	180
6.2	Les réponses modales	184
6.2.1	La sélection des réponses modales	184
6.2.2	Mise en oeuvre et exemples	186
6.2.3	Autres exemples	190
Chapitre 7 : Partitions longitudinales, contiguïté		197
7.1	Les trois structures de base	197
7.2	Homogénéité des valeurs d'une variable	199
7.2.1	Graphe associé à une structure de contiguïté	200
7.2.2	Matrices de contiguïté	200
7.2.3	Le coefficient de contiguïté	202
7.2.4	Moments du coefficient de contiguïté	204
7.2.5	Un cas particulier : les séries temporelles	204

7.2.6	Utilisation du coefficient c	205
7.3	Homogénéité des facteurs en fonction d'une structure a priori	205
7.3.1	Homogénéité d'un facteur	206
7.3.2	Homogénéité des k premiers facteurs	206
7.4	Les agrégats et l'analyse de la contiguïté	207
7.5	Partitions longitudinales d'un corpus	209
7.5.1	Exemple de partition longitudinale	210
7.5.2	Analyse de la gradation "classe d'âge"	211
7.5.3	Spécificités connexes	213
7.6	Séries textuelles chronologiques	217
7.6.1	La série chronologique <i>Discours</i>	218
7.6.2	Spécificités chronologiques	220
7.6.3	Les accroissements spécifiques	221
7.6.4	Etude parallèle sur un corpus lemmatisé	224
7.7	Recherches en homogénéité d'auteur	226
7.7.1	Le corpus informatisé	229
7.7.2	Validation des résultats	234
7.7.3	Fragments de n chapitres consécutifs	236
7.7.4	Classification des chapitres	238
Chapitre 8	Analyse discriminante textuelle	241
8.1	Deux grandes familles de problèmes	242
8.1.1	Discrimination à partir de la forme : la stylométrie	243
8.1.2	Discrimination globale	244
8.2	Les unités et indices de la stylométrie	245
8.2.1	"Mots outil", parties du discours	245
8.2.2	La richesse du vocabulaire	247
8.3	Modèles statistiques en stylométrie : un exemple	248
8.3.1	Modélisations de la gamme des fréquences	248
8.3.2	Le problème d'attribution	249
8.3.3	Un modèle non-paramétrique d'estimation	252
8.3.4	Autres approches du problème	254
8.4	Analyses discriminantes globales	256
8.4.1	Principe général	256
8.4.2	Unités pour la discrimination globale	258
8.4.3	Discrimination et réponses modales	259
8.4.4	Discrimination régularisée par analyse des correspondances préalable	261
8.4.5	Validation d'une discrimination	262

8.5 Discrimination globale et validation	263
8.5.1 L'exemple et le problème	263
8.5.2 Vocabulaire et analyse pour Tokyo	266
8.5.3 Réalité des configurations	272
8.5.4 Analyse discriminante et matrices de confusion	276
8.5.5 Conclusions	282
Annexe A Description sommaire de quatre logiciels	283
Annexe B Esquisse des algorithmes et structures de données	299
Glossaire	311
Références bibliographiques	321
Bibliographie complémentaire	332
Index des auteurs	335
Index des matières	339

Introduction

Les méthodes de *statistique textuelle* rassemblées dans le présent ouvrage sont nées de la rencontre entre plusieurs disciplines : l'étude des textes, la linguistique, l'analyse du discours, la statistique, l'informatique, le traitement des enquêtes, pour ne citer que les principales. Notre démarche s'appuie à la fois sur les travaux d'un courant aux dénominations changeantes (*statistique lexicale, statistique linguistique, linguistique quantitative, etc.*) qui associe depuis une cinquantaine d'années la méthode statistique à l'étude des textes, et sur l'un des courants de la statistique moderne, la *statistique multidimensionnelle*.

L'outil informatique est aujourd'hui utilisé par un nombre croissant d'utilisateurs pour des tâches qui impliquent la saisie et le traitement de grands ensembles de textes. Cette diffusion renforce à son tour la demande d'outils de gestion et d'analyse des textes qui émane des praticiens et des chercheurs de nombreuses disciplines. Confrontés à des textes nombreux recueillis dans des enquêtes socio-économiques, des entretiens, des investigations littéraires, des archives historiques ou des bases documentaires, ces derniers attendent en effet une aide en matière de classement, de description, de comparaisons...

Nous tenterons précisément de montrer comment les possibilités actuelles de calcul et de gestion peuvent aider à décrire, assimiler et enfin à critiquer l'information de type textuel.

Le choix d'une stratégie de recherche ne peut être opéré qu'en fonction d'objectifs bien définis. Quel type de texte analyse-t-on ? Pour tenter de répondre à quelles questions ? Désire-t-on étudier le vocabulaire d'un texte en vue d'en faire un commentaire stylistique ? Cherche-t-on à repérer des *contenus* à travers les réponses à un questionnaire ? S'agit-il de mettre en évidence les motivations pour l'achat d'un produit à partir d'opinions exprimées dans des entretiens ? Ou de classer des documents afin de mieux les retrouver ultérieurement ?

Bien entendu, aucune méthode d'analyse figée une fois pour toutes ne saurait répondre entièrement à des objectifs aussi diversifiés. Il nous est apparu

cependant qu'un même ensemble de méthodes apportait dans un grand nombre d'analyses de caractère textuel un éclairage irremplaçable pour avancer vers la solution des problèmes évoqués.

L'ouvrage que nous avons publié chez le même éditeur en 1988 sous le titre *Analyse statistique des données textuelles* concernait essentiellement l'analyse exploratoire des réponses aux questions ouvertes dans les enquêtes. Le contenu en a été élargi tant au niveau de la méthodologie qu'en ce qui concerne les domaines d'application.

Dans ce nouvel exposé, il ne s'agit plus uniquement de décrire et d'explorer, mais aussi de mettre à l'épreuve les hypothèses, de prouver la réalité de traits structuraux, de procéder à des prévisions. Quant au champ d'application des méthodes présentées, il dépasse dorénavant le cadre des traitements des réponses à des questions ouvertes et concerne des corpus de textes beaucoup plus généraux. Enfin, on a tenté de prendre en compte les travaux qui ont été réalisés depuis la parution du premier ouvrage.

L'accès à de nouveaux champs d'application, même lorsqu'il s'agit de méthodes éprouvées, peut demander une préparation des matériaux statistiques, un effort de clarification conceptuelle, une économie dans l'agencement des algorithmes, une sélection et une présentation spécifique des résultats. Ceci est tout particulièrement vrai pour ce qui concerne le domaine des études textuelles. Dans ce domaine en effet, la notion de *donnée* qui est à la base des comptages statistiques doit faire l'objet d'une réflexion spécifique.

D'une part il est nécessaire de découper des unités dans la chaîne textuelle pour réaliser des comptages utilisables par les analyses statistiques ultérieures. De l'autre, la chaîne textuelle ne peut être réduite à une succession d'unités n'ayant aucun lien les unes avec les autres car beaucoup des *effets de sens* du texte résultent justement de la disposition relative des formes, de leurs juxtapositions ou de leurs cooccurrences éventuelles.

* * *

Le premier chapitre, *Domaines et problèmes*, évoque à la fois : les domaines disciplinaires concernés (linguistique, statistique, informatique), les problèmes et les approches. Il précise dans chaque cas la nature du *matériau de base* que constituent les textes rassemblés en corpus.

Le second chapitre, *Les unités de la statistique textuelle*, est consacré à l'étude des unités statistiques que les programmes lexicométriques devront découper ou reconnaître (formes, segments répétés). Il aborde les aspects fondamentaux de l'approche quantitative des textes, les propriétés de ces unités ; il précise leurs pertinences respectives en fonction des champs d'application.

Les troisième et quatrième chapitres, *L'analyse des correspondances des tableaux lexicaux*, et *La classification automatique des formes et des textes*, présentent les techniques de base de l'*analyse statistique exploratoire* des données multidimensionnelles à partir d'exemples que l'on a souhaité les plus simples possibles.

Le cinquième chapitre : *Typologies, visualisations*, applique les outils présentés aux chapitres trois et quatre à la description des associations entre formes et entre catégories. Il fournit des exemples d'application *en vraie grandeur* commentés du point de vue de la méthode statistique. Il détaille les règles de lecture et d'interprétation des résultats obtenus, fait le point sur leur portée méthodologique.

Pour compléter ces représentations synthétiques, le sixième chapitre, *Éléments caractéristiques, réponses ou textes modaux*, présente les calculs dits de *spécificité* ou de *formes caractéristiques* qui permettent de repérer, pour chacune des parties d'un corpus, celles des unités qui se signalent par leurs fréquences atypiques. La sélection automatique des *réponses modales* ou des textes modaux permet de replacer les formes dans leur contexte, et de caractériser, lorsque cela est possible, des parties de texte, en général volumineuses, par des portions plus petites (phrases, paragraphes, documents, réponses dans le cas d'enquêtes). On résume ainsi, dans le cas des réponses libres, l'ensemble des réponses d'une catégorie de répondants par quelques réponses effectivement attestées dans le corpus, choisies en raison de leur caractère représentatif.

Le septième chapitre, *Partitions longitudinales, contiguïté*, traite le problème des informations *a priori* qui concernent les parties d'un corpus. Dans de nombreuses applications, en effet, l'analyste possède, avant toute démarche de type quantitatif, des informations qui lui permettent de rapprocher entre elles certaines des parties, ou encore de dégager un ordre privilégié parmi ces dernières (*séries textuelles chronologiques*). On étudie dans ce chapitre, en présentant une méthode et de nombreux exemples d'application, les relations de dépendance que l'on peut observer entre ces structures et les profils lexicaux des parties.

Enfin le huitième chapitre, consacré à l'*Analyse discriminante textuelle*, étudie, au sens statistique du terme, le *pouvoir de discrimination* des textes. Comment affecter un texte à un auteur (ou à une période) ? Peut-on prévoir l'appartenance d'un individu à une catégorie à partir de sa réponse à une question ouverte ? Comment classer (ici : affecter à des classes préexistantes) un document dans une base de données textuelles ? On tente dans ce chapitre, qui contient des exemples d'application variés, de montrer quels sont les apports de la statistique textuelle à la stylométrie, à la recherche documentaire, ainsi qu'à certains modèles prévisionnels.

Le cheminement méthodologique auquel nous invitons le lecteur verra ses étapes illustrées par des corpus de textes provenant de sphères de recherche très différentes. Les résultats présentés à ces occasions concernent des textes littéraires, des corpus de réponses libres dans des enquêtes françaises et internationales, des discours politiques.

L'ensemble des exemples devrait permettre au lecteur d'apprécier la variété des applications réalisées et potentielles, la complémentarité des divers traitements, tout en progressant dans l'assimilation et la maîtrise des méthodes, et surtout dans sa capacité à évaluer et critiquer les résultats.