



**HAL**  
open science

## La théorie de l'information vingt ans après Guiraud

Étienne Brunet

► **To cite this version:**

Étienne Brunet. La théorie de l'information vingt ans après Guiraud. Annales de la Faculté des Lettres et Sciences Humaines de Nice, 1985, Hommage à Pierre Guiraud, 52, pp.89-109. hal-01575406

**HAL Id: hal-01575406**

**<https://hal.science/hal-01575406>**

Submitted on 19 Aug 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.

## La théorie de l'information vingt ans après Guiraud

Etienne Brunet

Tous les raffinements mathématiques dont on peut entourer le rapport  $N/V$  ne peuvent dissoudre l'ambiguïté initiale qui s'attache à la notion de richesse lexicale<sup>1</sup>. Il ne suffit pas de réduire l'influence de la longueur des textes comparés. Encore faudrait-il s'attarder à la qualité, concrète ou abstraite, des vocables rencontrés et ne pas confondre l'abondance des objets et l'abondance des mots. En outre, en s'en tenant à la seule structure lexicale, un même rapport  $N/V$  peut être obtenu de bien des façons différentes, suivant qu'un texte privilégie les fréquences moyennes ou extrêmes. Enfin ce rapport peut varier dans le déroulement d'un texte et ce mouvement peut être intéressant à observer.

C'est pourquoi on peut penser que pour rendre compte de la diversité lexicale d'un texte il convient non seulement d'explorer toutes les fréquences – ce qu'on fait partiellement dans la loi binomiale et dans l'indice de Yule-Herdan –, mais aussi de suivre la progression lexicale de ce texte à des intervalles réguliers. Une recherche dans cette direction a été menée par Étienne Évrard sur un corpus d'auteurs latins et exposée au 11<sup>e</sup> Colloque de l'ALLC, à Louvain. Évrard a songé à utiliser ici la

---

<sup>1</sup> Ce rapport n'est jamais brut car il serait trop dépendant de l'étendue du texte (c'est-à-dire du nombre d'occurrences  $N$ , qui croît plus vite que celui des vocables  $V$ ). Diverses pondérations ont été apportées depuis que Guiraud a proposé la première :

$$V = \sqrt{N} = \text{constante } 22$$

Voir P. Guiraud, *Problèmes et méthodes de la statistique linguistique*, p. 89. Parmi ces indices de richesse lexicale, on citera, parmi bien d'autres, l'indice Uber de D. Dugast, celui de Rubet, celui de Yule-Herdan, et le nôtre ou indice  $w$ . Toutes les formules proposées sont, à l'image de celle de Guiraud, des approximations expérimentales. Une seule, prônée par Ch. Muller, procède par raisonnement et prend appui sur une loi générale (formule binomiale).

notion d'entropie<sup>2</sup> qui évoluerait entre deux limites : la valeur  $\log N$ , lorsque la diversité lexicale est à son maximum (lorsque toutes les occurrences appartiennent à des vocables différents) et la valeur zéro (lorsque le texte n'a qu'un seul vocable). De la forme initiale de l'entropie :

$$H = - \sum n_i p_i \log(p_i)$$

(où  $n$  désigne l'effectif d'une classe de fréquence  $x_i$ , et  $p_i$  la probabilité attachée à un mot de fréquence  $x_i$  :  $p_i = x_i / N$ ), Évrard aboutit à la formule :

$$H = \log(N) - 1/N \sum n_i x_i \log(x_i)$$

Prenons par exemple les 20 premières pages de la *Faute de l'abbé Mouret*. L'on y trouve 982 vocables de fréquence 1, 285 pour  $f = 2$ , 114 pour  $f = 3$ , 76 pour  $f = 4$ , etc. Le cumul de toutes les classes jusqu'à la plus élevée ( $f = 74$ ) aboutit à un total de 4 309 occurrences dont le logarithme (décimal) est de 3,6344. Le terme  $n_i x_i \log(x_i)$  vaut : 2 789,22 et pondéré par  $N$  : 0,6473. L'entropie a donc pour valeur :

$$H = 3,6344 - 0,6473 = 2,9871$$

On a procédé à de tels calculs non seulement pour les 20 premières pages du texte en question, mais aussi pour le roman entier, saisi à intervalles réguliers, de 10 en 10 pages, et même pour tous les textes qui constituent le cycle des *Rougon-Macquart* (en y adjoignant deux romans antérieurs : *Thérèse Raquin* et *Madeleine Férat*). Comme on ne saurait restituer ici les 200 000 cases du tableau des effectifs, on s'est contenté de reproduire la progression de  $N$  dans le tableau 1 et celle de  $V$  dans le tableau 2. Avant d'aborder les résultats consignés dans le tableau 5, il convient d'avertir le lecteur que les calculs sont considérables puisqu'à chaque pas, pour chaque texte, il faut réordonner tout le tableau de distribution des fréquences. Imaginons le discours comme un sablier dont le contenu se répand. La présente expérimentation consiste à boucher le sas à intervalles réguliers et à examiner la position de chaque grain de sable dans la pyramide qui s'est constituée et dont la base représente la donnée  $V$ , le volume la donnée  $N$  et la hauteur la fréquence la plus élevée. Un même volume peut être obtenu avec une base plus ou moins large et

---

<sup>2</sup> Là encore Guiraud a été le premier linguiste en France à utiliser la notion d'entropie et à mettre en œuvre la théorie de l'information. Voir *Problèmes...* p. 79.

une hauteur plus ou moins grande. Chaque colonne de grains empilés représente une classe de fréquence, les hapax étant les grains libres qui sont répandus à la périphérie et qui n'en supportent (provisoirement) aucun autre. Or il existe bien des formes de pyramides, de la plus plate à la plus aiguë, et l'entropie mesure précisément la forme de cette pyramide et principalement le rapport entre la surface de base et la hauteur.

Mais le sable est mouvant et la pyramide changeante. Généralement aplatie au départ, elle prend ensuite une forme plus élancée et Étienne Évrard tente d'imposer une mesure à cette évolution. Avant de suivre ses pas, nous avons pris quelques précautions du côté des hautes fréquences dont l'influence est prépondérante sur l'entropie comme elle l'est sur l'indice de Yule-Herdan. Observons en effet ce qui se passe dans le cas du *Docteur Pascal* (tableau 3). La sommation  $S = \sum n_i x_i \cdot \log(x_i)$  atteint un maximum de 105 002 quand le texte touche à sa fin, à la 310<sup>e</sup> page, et que le total des occurrences est de 65 620. Si l'on prête attention, on remarque qu'il manque à l'appel 54 391 occurrences, puisque l'étendue de ce texte est en réalité de 120 011. C'est que nous avons écarté volontairement une trentaine de mots très fréquents dont la masse eût troublé les mesures de l'entropie<sup>3</sup>. Prenons l'exemple de la préposition DE qui a 6 663 emplois dans le *Docteur Pascal*. La part de ce seul mot dans la sommation eût été écrasante soit :

$$6\,663 \times \log 6\,663 = 25\,477$$

c'est-à-dire le quart du total obtenu. Même en écartant ces mots encombrants, le calcul reste sensible aux mots fréquents. Ceux qui ont plus de 100 occurrences (il y en a 111 dans le texte considéré, sans compter les 36 unités écartées) représentent la moitié du total :  $53\,086/105\,002 = 0,51^4$ .

---

<sup>3</sup> Voici la liste des mots écartés : A, AVOIR, CE, DANS, DE, DES (et DES), DU (et DU), ELLE, EN, ET, ETRE, IL, ILS, JE, L', LA (et LA), LE, LES, LUI, NE (et NE), PAS, PLUS, POUR, QUE, QUI, SA, SE, SES, SON. UN, NE, VOUS. L'ensemble de ces 36 vocables représente 44% des occurrences du corpus.

<sup>4</sup> Bien entendu les fréquences supérieures à 100 entrent dans le calcul avec leurs valeurs et leurs effectifs propres. Leur influence sur le terme  $n_i x_i \cdot \log(x_i)$  croît avec l'étendue :

Pages	10	50	100	150	200	250	310
Influence	0%	2%	27%	33%	40%	42%	51%

Tableau 1 : Les occurrences de 10 en 10 pages

Table with columns labeled RAQ, PER, FOR, CUR, VEN, CON, FAU, EXC, ASS, FAG, NAN, POT, BON, JOI, GER, OEU, TER, REV, BET, ARG, DEB, PAS and rows of numerical data.

Tableau 2 : Les vocables de 10 en 10 pages

Table with columns labeled RAQ, PER, FOR, CUR, VEN, CON, FAU, EXC, ASS, FAG, NAN, POT, BON, JOI, GER, OEU, TER, REV, BET, ARG, DEB, PAS and rows of numerical data.

Tableau 3 : L'entropie. Détail du calcul pour le *Docteur Pascal*

	N	v	$v/\sqrt{N}$	H réel	H théo	$\log(N)$	S	S/N
1	2141	1061	22.9301	2.8577	2.8607	3.3306	1012	0.4729
2	4322	1605	24.4136	2.9458	2.9807	3.6357	2981	0.6899
3	6529	2179	26.9671	3.0446	3.0401	3.8148	5028	0.7702
4	8737	2573	27.5270	3.0815	3.0770	3.9414	7512	0.8599
5	10811	2924	28.1219	3.1132	3.1015	4.0339	9953	0.9207
6	12911	3238	28.4968	3.1308	3.1202	4.1110	12654	0.9801
7	14962	3492	28.5482	3.1404	3.1346	4.1750	15479	1.0346
8	17206	3700	28.2073	3.1444	3.1472	4.2357	18777	1.0913
9	19375	3907	28.0687	3.1495	3.1573	4.2872	22043	1.1378
10	21588	4263	29.0141	3.1799	3.1658	4.3342	24919	1.1543
11	23799	4449	28.8392	3.1852	3.1730	4.3766	28354	1.1914
12	25975	4612	28.6162	3.1909	3.1790	4.4146	31784	1.2237
13	28007	4718	28.1919	3.1904	3.1840	4.4473	35201	1.2569
14	30212	4871	28.0239	3.1955	3.1887	4.4802	38812	1.2847
15	32330	4976	27.6743	3.1926	3.1926	4.5096	42577	1.3170
16	34586	5120	27.5309	3.1959	3.1963	4.5389	46449	1.3430
17	36776	5238	27.3139	3.1979	3.1995	4.5656	50298	1.3677
18	38924	5368	27.2084	3.2040	3.2023	4.5902	53957	1.3862
19	41105	5477	27.0144	3.2082	3.2048	4.6139	57780	1.4057
20	43203	5596	26.9228	3.2101	3.2070	4.6355	61581	1.4254
21	45377	5696	26.7394	3.2103	3.2091	4.6568	65641	1.4466
22	47582	5786	26.5251	3.2096	3.2109	4.6774	69843	1.4679
23	49763	5847	26.2108	3.2061	3.2126	4.6969	74186	1.4908
24	51890	5923	26.0016	3.2055	3.2140	4.7151	78334	1.5096
25	54040	5993	25.7802	3.2077	3.2154	4.7327	82411	1.5250
26	56203	6078	25.6378	3.2075	3.2166	4.7498	86677	1.5422
27	58297	6162	25.5210	3.2106	3.2176	4.7656	90651	1.5550
28	60400	6229	25.3454	3.2104	3.2186	4.7810	94866	1.5706
29	62553	6321	25.2733	3.2126	3.2195	4.7962	99059	1.5836
30	64827	6404	25.1520	3.2159	3.2204	4.8118	103452	1.5958
31	65620	6433	25.1128	3.2169	3.2206	4.8170	105002	1.6002

$S = \sum_{i=1}^n x_i \log(x_i)$       pour  $a = -0.02415$  et  $b = 0.55823$   
 $r = 0.9935$        $H \text{ max} = -b^2 / 4a = 3.2259$   
 $x \text{ max} = -b/2a = 11.5576$

TABLEAU DE DISTRIBUTION pour le texte entier

(pour f= 1 à 100)

2210	991	642	416	282	236	175	147	125	103	76	84	68	48	44	43	31	36	30	29
31	26	30	20	18	19	24	7	10	9	11	12	18	9	9	10	11	9	10	9
13	8	8	5	12	6	5	5	3	5	6	3	5	4	4	6	2	4	1	3
5	7	4	3	5	6	2	2	4	5	3	1	3	2	4	3	2	1	0	4
2	1	4	0	1	4	1	1	1	2	1	0	1	3	1	1	1	2	3	111

Quant à l'unité de mesure choisie – tranches de 10 pages – elle est arbitraire, comme l'eût été l'adoption des lignes, des mots ou des lettres. Mais comme il s'agit d'une édition homogène (celle de la Pléiade), la mesure est à peu près constante<sup>5</sup> et le choix imposé par des contraintes

<sup>5</sup> A deux réserves près : d'une part les deux premiers textes (*Thérèse Raquin* et *Madeleine Féral*) sont empruntés à une autre édition, et d'autre part là où le dialogue est

techniques<sup>6</sup> en vaut un autre. De toute façon le nombre d'occurrences a été relevé dans chaque tranche pour servir de pondération lorsqu'il en serait besoin.

On pourrait s'attarder sur le tableau de distribution à chaque palier de l'observation. Rarement pareille occasion a été donnée d'étudier dans le détail et au ralenti le mouvement des fréquences au fur et à mesure qu'un discours s'organise dans la durée. Mais l'information est si abondante qu'on ne saurait en rendre compte sans faire appel à des programmes sophistiqués de traitement graphique. Seule une image animée pourrait montrer en effet la forme changeante de la pyramide des mots et les modifications incessantes – amoncellements, excroissances, tassements, écroulements – que la consistance variable de la matière lexicale peut engendrer à tel ou tel point de la surface. Des ondes se créent, l'effectif d'une classe se gonflant quand celui de la classe inférieure s'épuise, en attendant qu'un afflux nouveau comble les vides. Cette mécanique ondulatoire explique les variations d'effectifs qui affectent chaque classe de fréquence et qui créent parfois des baisses provisoires dans un mouvement général ascendant. Voici à titre d'exemple l'évolution des effectifs de la classe 20 dans la *Faute de l'abbé Mouret* :

0,0,2,1,6,9,8,9,13,10,8,13,12,19,15,17,14,15,14,17,16,21,18,17,25,25,36,30,29,31,27,29,29<sup>7</sup>.

La seule classe qui paraît échapper à la baisse est celle de la fréquence 1. Dans la *Faute de l'abbé Mouret*, l'effectif des hapax croît pratiquement sans cesse de 665 pour la première tranche à 2 284 pour l'ensemble du texte. Mais il y a des à-coups, des accélérations et des ralentissements de la progression de cette classe qui est très sensible au renouvellement lexical. Ainsi, dans le même roman, l'épisode du *Paradou*, si riche en couleurs, en odeurs et en mots, produit un accès de

---

abondant, les alinéas se multiplient – à chaque réplique – et la densité en mots de chaque page s'affaiblit. Voir dans le tableau 7 le relevé des occurrences et dans le tableau 8 le relevé des vocables.

<sup>6</sup> Le programme prend ses données dans l'index où seule subsiste l'indication de la page mais non le numéro d'ordre de chaque mot.

<sup>7</sup> Pour saisir en mouvement les phénomènes de transit qui se produisent dans les classes de fréquence, on s'est intéressé, à titre d'exemple, aux gains et aux pertes de la classe 10 au fur et à mesure que se développent les 22 textes du corpus. Il s'agit certes d'un mouvement ascendant, les effectifs se gonflant de la première à la dernière tranche. Mais c'est aussi un mouvement oscillatoire ou ondulatoire, fait de gains provisoires et de décrues passagères. La progression est hésitante et vibratile avec deux pas en avant pour un en arrière (593 écarts positifs pour 225 négatifs). Les classes de fréquence sont des lieux de passage, des vases communicants où se déploie la mécanique des fluides.







	H	rang 1	rang 2		H	rang 1	rang 2
RAQ	3,168	20	19	POT	3,175	19	20
FER	3,177	18	22	BON	3,238	9	10
FOR	3,258	4	6	JOI	3,198	16	15
CUR	3,260	3	2	GER	3,255	5	9
VEN	3,243	7	3	OEU	3,275	1	1
CON	3,138	22	21	TER	3,230	10	8
FAU	3,200	15	12	REV	3,221	11	7
EXC	3,249	6	5	BET	3,205	14	17
ASS	3,194	17	14	ARG	3,265	2	4
PAG	3,162	21	18	DEB	3,243	8	16
NAN	3,207	13	11	PAS	3,217	12	13

On a rapproché le classement obtenu selon l'entropie (rang 1) de celui que la loi binomiale nous avait proposé pour la richesse lexicale. La similitude est frappante et l'entropie semble donc bien mesurer, comme le pense Évrard, la diversité lexicale d'un texte. Il y a pourtant quelques distorsions et le problème est de savoir à laquelle des deux méthodes attribuer la déviance. Il ne sert à rien de solliciter le témoignage des autres indices (W, Uber ou Rubet) puisque, s'appuyant sur  $N$  et  $V$  comme la loi binomiale, ils ne peuvent que souscrire au classement de cette dernière. Tous ces indices luttent à l'envi contre l'influence de la longueur et sur ce point l'entropie n'est pas sans défaut. Observons en effet la progression de  $H$  sur une même colonne, c'est-à-dire dans un même texte. Partout l'entropie est croissante d'un mouvement régulier<sup>8</sup> qui ne souffre aucune exception. Et cela tient à sa définition, l'entropie partant nécessairement de zéro (lorsque le texte n'a qu'une occurrence et qu'un vocable) et s'élevant jusqu'à un maximum de  $\log N$ , qu'elle est loin d'atteindre (dans l'ordre de grandeur de nos textes, l'entropie se situe à 3 environ, alors que la limite est fixée près de 4). Et cette caractéristique nuit à l'utilisation de l'entropie quand les textes comparés sont de longueur différente. Ainsi, à son terme, la *Débâcle*, parvient à une entropie de 3,243 qui dépasse celle du *Rêve* (3,221) dont le vocabulaire est plus riche mais dont l'étendue est trois fois moindre. Mais si dans la *Débâcle*, on prend la valeur de  $H$  au niveau du *Rêve*, c'est-à-dire au bout de 200 pages, l'entropie n'est plus que de 3,186. Cette influence de la longueur fait ainsi gagner 8 places à la *Débâcle* dans le classement de l'entropie, et 4 à *Germinal*, et dans le même temps elle pénalise le *Rêve* de 4 places.

---

<sup>8</sup> Il ne s'agit pas d'un mouvement uniforme mais d'un mouvement uniformément ralenti – ce qui invite à en rendre compte par une équation du second degré.

Tableau 6 : Classements selon l'entropie

PAGE	HAQ	FER	POR	CUR	VEN	CON	PAU	EXC	ASS	PAG	NAN	POT	BON	JOI	GER	OEU	TER	REV	BET	ARG	DEB	PAS
10	19	14	3	9	18	22	10	15	13	21	12	20	11	6	17	1	16	2	5	4	8	7
20	19	20	7	3	18	22	5	11	17	21	13	16	12	9	8	2	10	1	14	4	6	15
30	19	14	8	3	18	22	9	11	20	21	10	17	15	13	7	2	16	1	12	4	6	5
40	19	15	7	4	12	22	14	11	20	21	16	17	13	9	10	2	8	1	18	3	6	5
50	20	11	5	3	10	22	15	14	19	21	17	18	12	9	13	4	8	1	16	2	6	7
60	21	12	5	3	10	22	15	14	18	20	17	19	13	8	11	4	9	1	16	2	6	7
70	20	14	5	3	10	22	16	12	17	21	19	18	13	8	11	4	6	1	15	2	7	9
80	18	15	4	1	10	22	13	12	16	21	20	19	14	7	9	5	6	2	17	3	8	11
90	18	16	4	1	9	22	8	11	17	21	20	19	14	7	12	5	6	3	15	2	10	13
100	18	15	3	1	7	22	10	12	17	20	21	19	13	9	14	5	6	4	16	2	11	8
110	17	15	3	1	7	22	10	11	18	19	21	20	13	9	14	5	6	4	16	2	12	8
120	17	13	2	1	6	22	10	11	18	19	20	21	12	9	14	4	7	5	16	3	15	8
130	17	14	2	1	6	22	11	9	18	19	20	21	12	10	13	4	7	5	16	3	15	8
140	18	16	2	1	7	22	9	10	17	20	19	21	12	11	13	4	6	5	14	3	15	9
150	18	16	2	1	6	22	8	10	17	21	19	20	11	12	13	4	7	5	14	3	15	9
160	18	16	3	1	7	22	8	9	17	21	19	20	11	13	12	4	5	6	14	2	15	10
170	19	16	3	1	7	22	10	8	17	21	18	20	11	14	12	4	5	6	13	2	15	9
180	18	16	3	1	5	22	10	8	19	20	17	21	12	14	11	4	7	6	13	2	15	9
190	19	16	3	1	5	22	11	9	18	21	17	20	12	14	10	4	7	6	13	2	15	8
200	18	17	2	1	5	22	12	9	19	20	16	21	11	13	10	4	7	6	14	3	15	8
210	18	17	3	1	5	22	12	8	19	21	16	20	11	13	10	4	7	6	14	2	15	9
220	18	19	3	2	5	22	13	9	17	21	16	20	11	12	8	4	6	7	15	1	14	10
230	17	19	3	1	5	22	12	9	18	21	16	20	11	13	8	4	6	7	15	2	14	10
240	18	19	4	1	5	22	12	9	17	21	16	20	10	14	8	2	6	7	15	3	13	11
250	19	17	3	1	5	22	13	9	18	21	16	20	10	14	6	4	7	8	15	2	12	11
260	19	17	2	1	5	22	15	8	18	21	16	20	10	12	6	3	7	8	14	4	13	11
270	19	17	3	1	5	22	15	8	18	21	16	20	11	12	6	2	7	9	13	4	14	10
280	19	17	3	2	5	22	14	7	18	21	16	20	10	15	6	1	8	9	13	4	12	11
290	19	17	3	2	5	22	15	6	18	21	16	20	10	13	7	1	8	9	14	4	12	11
300	19	17	3	2	5	22	14	7	18	21	16	20	11	13	6	1	8	9	15	4	12	10
310	19	17	3	2	5	22	13	7	18	21	16	20	11	14	6	1	8	9	15	4	12	10
320	19	17	3	2	5	22	13	7	18	21	16	20	11	15	6	1	8	9	14	4	12	10
330	19	18	4	2	5	22	14	7	17	21	16	20	10	15	6	1	8	9	13	3	12	11
340	19	18	3	2	5	22	14	7	17	21	15	20	9	16	7	1	8	10	13	3	12	11
350	19	18	4	2	6	22	14	5	17	21	15	20	9	16	7	1	8	10	13	2	11	12
360	19	18	4	3	6	22	15	5	17	21	14	20	9	16	7	1	8	10	13	2	11	12
370	20	18	4	3	7	22	15	5	17	21	13	19	9	16	6	1	8	10	14	2	11	12
380	20	18	4	3	7	22	15	5	17	21	13	19	8	16	6	1	9	10	14	2	11	12
390	20	18	4	3	7	22	15	5	17	21	13	19	8	16	6	1	9	10	14	2	11	12
400	20	18	4	3	7	22	15	5	17	21	13	19	8	16	6	1	9	10	14	2	11	12
410	20	18	4	3	7	22	15	5	17	21	13	19	8	16	6	1	9	11	14	2	10	12
420	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	9	11	14	2	10	12
430	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	9	11	14	2	10	12
440	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	9	11	14	2	10	12
450	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	9	11	14	2	10	12
460	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	10	11	14	2	9	12
470	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	10	11	14	2	9	12
480	20	18	4	3	7	22	15	6	17	21	13	19	8	16	5	1	10	11	14	2	9	12
490	20	18	4	3	7	22	15	6	17	21	13	19	9	16	5	1	10	11	14	2	8	12
500	20	18	4	3	7	22	15	6	17	21	13	19	9	16	5	1	10	11	14	2	8	12
510	20	18	4	3	7	22	15	6	17	21	13	19	9	16	5	1	10	11	14	2	8	12
520	20	18	4	3	7	22	15	6	17	21	13	19	9	16	5	1	10	11	14	2	8	12

C'est pourquoi les classements n'ont de valeur qu'au même niveau et nous les avons calculés tranche après tranche dans le tableau 6. On les lira sur les lignes de cette matrice. La première montre que dès la dixième page l'*Œuvre* a pris la tête du corpus au regard de la diversité lexicale. Elle la cède ensuite au *Rêve*, puis à la *Curée* avant de la retrouver définitivement. En queue de classement au contraire, le même texte (la *Conquête*) se maintient d'un bout à l'autre. Si l'on veut à toute force fixer un classement, on peut choisir un niveau qui représente le meilleur compromis entre les textes longs et courts. Ce pourrait être par exemple la 30<sup>e</sup> ligne qui indique pour chaque texte l'entropie atteinte à la page 300. On se rapproche un peu alors du classement de la loi binomiale

(0,91). Mais le tableau 6 est aussi intéressant par ses colonnes qui montrent l'évolution interne de chaque texte. Certains, mal classés au départ, améliorent leur position au fil des pages. Certains autres au contraire déçoivent, sans compter et sans attendre, leurs richesses lexicales et s'épuisent par la suite. Le *Ventre de Paris* a un profil caractéristique du premier groupe, comme aussi *Germinal*, alors que le *Rêve* représente le second.

Ce qui sépare l'entropie des autres indices tient aussi au fait que la mesure ne porte pas exactement sur le même objet. Tout indice qui met en rapport  $N$  et  $V$  est nécessairement très sensible aux basses fréquences qui influent grandement sur  $V$  (les seuls hapax, suivant les cas, représentent 30 à 50 pour cent des vocables). Le calcul de l'entropie fait au contraire la part belle aux hautes fréquences, au point que nous avons dû filtrer les fréquences extrêmes. Il tend à négliger par contre les basses fréquences et en particulier se montre totalement indifférent aux hapax car le terme  $n_i x_i \log(x_i)$  est nul lorsque  $x_i$  vaut 1 ( $\log(1)$  valant zéro). Les textes où les hapax foisonnent, principalement l'*Assommoir* et le *Ventre de Paris* reculent donc dans le classement de l'entropie.

Mais Évrard ne destine pas l'entropie à ce rôle de classement. Car ce n'est pas pour lui la valeur finale de l'entropie qui importe (puisque'elle dépend de la longueur) mais bien plutôt le profil de cette entropie, la courbe qui rend compte d'une progression plus ou moins accentuée. Une fois de plus, comme dans le calcul de l'indice Uber, nous sommes ramenés à un problème d'ajustement. Et comme Dugast, Évrard songe à une fonction du second degré du type  $y = ax^2 + bx$  (il n'y a pas de terme  $c$  parce que la courbe commence nécessairement à 0, quand le texte n'a qu'un mot). Comme dans le cas de l'indice Uber, on aurait pu recourir à un autre modèle d'ajustement, par exemple à la fonction logarithmique  $y = a + b \log(x)$ ,  $y$  représentant l'entropie et  $x$  l'étendue du texte. Les résultats eussent été acceptables avec un coefficient de détermination proche de 0,90 (soit une corrélation de 0,95). Voir la partie supérieure du tableau 7. L'ajustement de la fonction du second degré est meilleur encore et les valeurs des paramètres  $a$  et  $b$  s'accordent avec les observations faites par Évrard sur son corpus latin :  $a$  est négatif (et  $b$  positif) et la concavité de la parabole est tournée vers le bas. Les valeurs relevées dans les 22 textes sont du même ordre et autorisent, semble-t-il, des classements auxquels Évrard donne une signification précise. Le paramètre  $b$  rendrait compte de la volonté plus ou moins accusée de diversification lexicale, alors que le paramètre  $a$  marquerait la limite de l'« essoufflement », la pesanteur qui infléchit la courbe vers le bas.

Tableau 7 : Entropie

APPROXIMATION LOGARITHMIQUE

		b	a	r <sup>2</sup>	r	rang b	rang a
1	RAQ	0.10637	2.07363	0.94242	0.97079	3	3
2	FER	0.08703	2.25783	0.88180	0.93904	21	18
3	FOR	0.10810	2.09573	0.91887	0.95858	2	4
4	CUR	0.10293	2.16360	0.85874	0.92668	7	10
5	VEN	0.11852	1.96605	0.94410	0.97165	1	1
6	CON	0.09995	2.03671	0.97900	0.98944	9	2
7	FAU	0.09260	2.20938	0.89745	0.94734	15	13
8	EXC	0.10487	2.09704	0.94706	0.97317	5	5
9	ASS	0.09224	2.18196	0.93381	0.96634	16	12
10	PAG	0.09824	2.11220	0.91461	0.95635	10	8
11	NAN	0.09032	2.20949	0.94676	0.97302	17	14
12	POT	0.08883	2.21502	0.90481	0.95121	19	16
13	BON	0.10018	2.13530	0.94659	0.97293	8	9
14	JOI	0.08771	2.26430	0.87405	0.93491	20	19
15	GER	0.10411	2.10131	0.94568	0.97246	6	6
16	OEU	0.08897	2.29717	0.91455	0.95632	18	21
17	TER	0.10492	2.10743	0.88828	0.94249	4	7
18	REV	0.09302	2.26865	0.79653	0.89249	14	20
19	BET	0.09593	2.16477	0.95784	0.97869	12	11
20	ARG	0.09652	2.22369	0.86762	0.93146	11	17
21	DEB	0.08367	2.30317	0.88176	0.93902	22	22
22	PAS	0.09361	2.21120	0.89373	0.94538	13	15

$y = a + b \log(x)$  pour  $y =$  entropie et  $x = N$

$r = 0,67$   
 $r = 0,88$   
 $r = 0,86$

CALCUL DES PARAMETRES A ET B

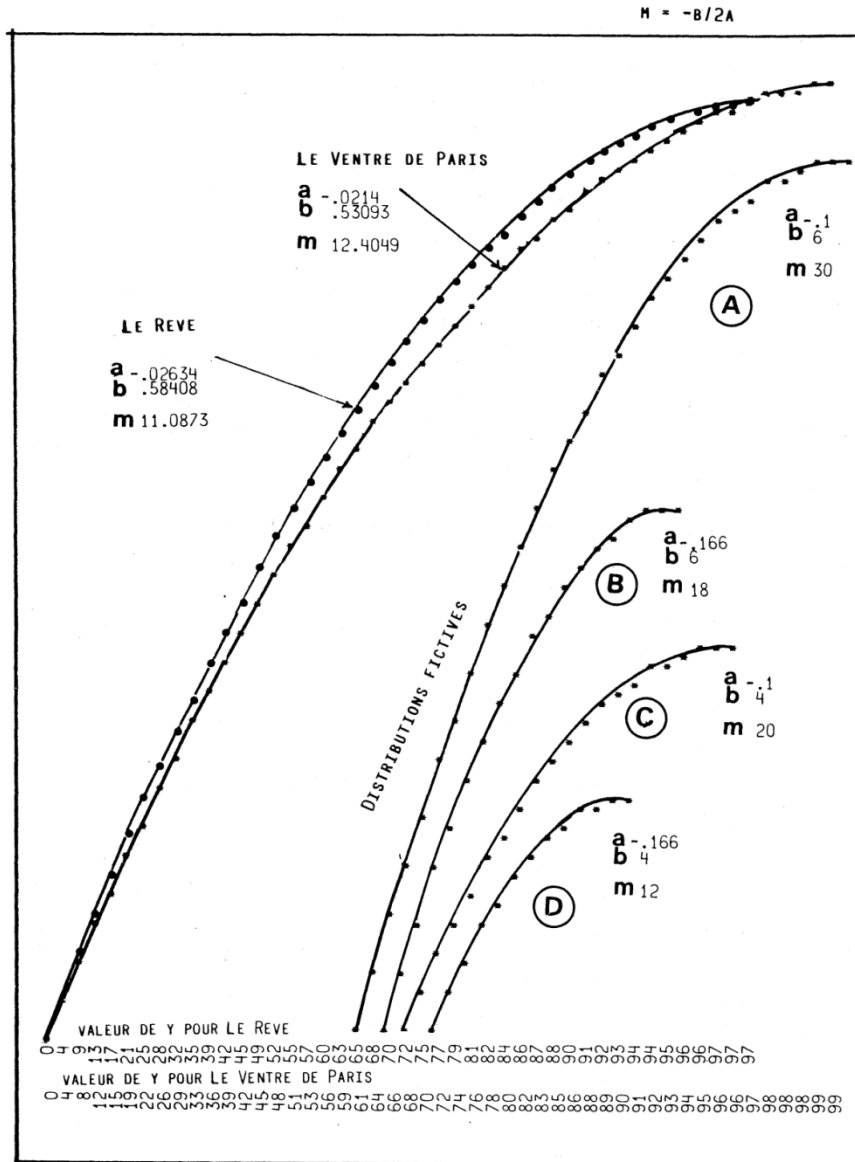
	N	log(N)	x <sup>max</sup>	rang x max	b	a	rang b	rang a
RAQ	37866	10.5418	11.3528	18	0.56151	-0.02473	16	17
FER	55619	10.9263	11.1375	21	0.57113	-0.02564	20	21
FOR	67902	11.1258	11.9886	2	0.54812	-0.02286	5	3
CUR	60778	11.0150	11.8222	5	0.55659	-0.02354	13	9
VEN	63921	11.0654	12.4049	1	0.53093	-0.02140	1	1
CON	64869	11.0801	11.7725	7	0.53094	-0.02255	2	2
FAU	68551	11.1353	11.5519	13	0.55634	-0.02408	11	12
EXC	73468	11.2046	11.8124	6	0.54857	-0.02322	6	6
ASS	92591	11.4359	11.5190	14	0.55222	-0.02397	10	11
PAG	59546	10.9945	11.6419	10	0.54554	-0.02343	3	7
NAN	82310	11.3182	11.4343	15	0.55662	-0.02434	14	14
POT	78642	11.2727	11.3903	16	0.55653	-0.02443	12	16
BON	86645	11.3696	11.7281	8	0.55005	-0.02345	9	8
JOI	68722	11.1378	11.3525	19	0.56490	-0.02488	18	18
GER	96246	11.4747	11.8460	4	0.54681	-0.02308	4	4
OEU	77252	11.2548	11.3883	17	0.57192	-0.02511	21	19
TER	94299	11.4542	11.8676	3	0.54876	-0.02312	7	5
REV	38554	10.5598	11.0873	22	0.58408	-0.02634	22	22
BET	73002	11.1982	11.6680	9	0.55003	-0.02357	8	10
ARG	83015	11.3268	11.5756	11	0.56489	-0.02440	17	15
DEB	108657	11.5960	11.2305	20	0.56961	-0.02536	19	20
PAS	65620	11.0916	11.5576	12	0.55823	-0.02415	15	13

$y = ax^2 + bx$   
 pour  $y =$  entropie et  $x = \log(N)$

$r = 0,96$

En théorie, les paramètres  $a$  et  $b$  ont effectivement cet effet sur une fonction du second degré. Et on le comprendra avec 4 exemples fictifs que nous avons projetés à droite de la figure 8. La pente initiale est surtout déterminée par le facteur  $b$  et elle est plus forte en  $A$  et  $B$  ( $b = 6$ ) qu'en  $C$  et  $D$  ( $b = 4$ ). La courbure, au contraire, est d'autant plus accusée que  $a$  est plus grand, et c'est le cas en  $B$  et  $D$  (où  $a$  vaut -0,166, contre -0,1 dans  $A$  et  $C$ ). La diversité maximale du vocabulaire suppose donc une pente raide ( $b$  fort) et une courbure peu prononcée ( $a$  faible). Une entropie restreinte (c'est-à-dire la sobriété lexicale) devrait au contraire gonfler le paramètre  $a$  et diminuer  $b$ . Or qu'observe-t-on ? Des situations intermédiaires du type  $B$  ou  $C$ , où  $a$  et  $b$  sont à la fois faibles ou forts. La partie gauche de la figure 8 représente les deux textes qui s'opposent le plus par la courbe de leur entropie : le *Rêve* et le *Ventre de Paris*. Le *Rêve* a les plus fortes valeurs de la série pour les deux paramètres  $a$  et  $b$  et le *Ventre de Paris* les deux plus faibles. Les deux courbes se différencient nettement sur le graphique : celle du *Rêve* s'élève et retombe plus rapidement, en une trajectoire plus verticale que la pesanteur épuise plus vite. Celle du *Ventre de Paris* est plus tendue et elle va plus loin. Or on observe les mêmes effets balistiques dans les 22 textes de notre corpus. Partout une relation constante est établie entre les facteurs  $b$  et  $a$ , entre l'angle d'attaque et l'épuisement de l'énergie. Si on fait un classement sur  $a$  et sur  $b$ , la corrélation est extrêmement forte ( $r = 0,96$ ). Dès lors il devient difficile de donner un sens particulier et différent à l'un et à l'autre paramètres, puisqu'ils vont de pair et qu'apparemment ils se rattachent à une même cause et rendent compte du même phénomène. D'ailleurs dans les données étudiées par Évrard, la liaison entre  $a$  et  $b$  paraît également établie. Précisons d'ailleurs que les paramètres  $a$  et  $b$  qu'on obtient par une approximation logarithmique sont eux aussi corrélés entre eux et qu'ils entretiennent des liens étroits avec ceux de la fonction du second degré. Dans le premier cas ( $a$  et  $b$ ) le coefficient est de 0,88, dans le second ( $a_1$  et  $a_2$ ,  $b_1$  et  $b_2$ ) la corrélation s'élève à 0,86 et 0,67 respectivement.

Tableau 8 : Courbe de l'entropie ( $y = ax^2 + bx$ )



Que conclure de notre coûteuse expérience ? Qu'il faut sans doute renoncer à distinguer les paramètres  $a$  et  $b$  puisque dans les faits observés ils donnent le même enseignement. Mais cette information n'est pas indifférente : la forme de la courbe de l'entropie renseigne sur la façon

dont un auteur gère le stock lexical, en prodigue ou en économe, en visant le court terme ou les longues échéances. Et de ce point de vue, les deux textes représentés dans le graphique 8, qui sont tous les deux parmi les plus riches des *Rougon-Macquart*, n'usent pas semblablement de leurs richesses. Dans l'univers de sacristie du *Rêve*, on sent que Zola est entré un peu par effraction, armé de dictionnaires et d'ouvrages spécialisés<sup>9</sup>. Il est pressé d'écouler son butin, presque à la sauvette, au lieu que dans le *Ventre de Paris*, Zola se sent chez lui, dans un monde familier dont il connaît depuis longtemps les inépuisables ressources. Il fait faire à son lecteur le tour du propriétaire mais il garde des réserves, des secrets et l'on sent qu'il pourrait offrir une seconde visite totalement différente de la première<sup>10</sup>.

Il existe cependant un moyen d'amplifier les minces différences qui peuvent séparer *a* et *b* afin de souligner le dessin de la courbe. Il s'agit tout simplement de projeter le point ultime que la courbe peut atteindre sur l'axe des *x* avant de rebrousser chemin, c'est-à-dire avant que la seconde branche de la parabole n'entame son mouvement de chute<sup>11</sup>.

---

<sup>9</sup> Le dossier préparatoire du *Rêve* est détaillé par H. Mitterand, pp. 1610-1615, la Pléiade tome IV. On y trouve un bric-à-brac de lectures sur la *Légende dorée*, l'architecture des cathédrales, les cérémonies religieuses, les règles administratives de l'adoption, les armoiries les vitraux, les broderies. Beaucoup de dictionnaires et d'encyclopédies dans cette documentation : l'*Art du brodeur* de Saint Aubain, l'*Almanach du Commerce*, le *Dictionnaire* de Savary, le *Dictionnaire de l'Industrie*, la *Légende dorée* de Jacques de Voragine, le *Grand Dictionnaire* de Pierre Larousse, le *Dictionnaire d'Architecture*, le *Dictionnaire des cérémonies et des rites sacrés*.

<sup>10</sup> La préparation livresque du *Ventre de Paris* est beaucoup plus légère : quelques ouvrages d'histoire sur le Second Empire, quelques notices sur les déportations en Guyane et quelques rapport sur l'administration des Halles, voir H. Mitterand tome 1, pp. 1622-1623. Les descriptions du *Ventre de Paris* sont faites *de visu* comme en témoigne Paul Alexis qui a souvent accompagné Zola dans la visite des Halles : « Un crayon à la main. Zola venait par tous les temps, par la pluie, le soleil, le brouillard, la neige, et à toutes les heures, le matin, l'après-midi, le soir, afin de noter les différents aspects. Puis, une fois, il y passa la nuit entière, pour assister au grand arrivage de la nourriture de Paris », cité par H. Mitterand, tome 1, p. 1617.

<sup>11</sup> Cela ne s'observe pas dans la réalité où l'entropie croît sans cesse, du moins dans l'ordre de grandeur où nous avons placé nos observations. C'est donc la portion ascendante de la courbe qui est utile, mais presque toute la portion ascendante, car on s'approche très près du sommet. Ainsi dans le *Rêve*, la limite  $\log N$  a pour valeur 10,5598 alors que le sommet théorique de la courbe est de :  $-b / 2a = 11,0873$ . Dans le cas du texte le plus long, la *Débâcle*, il semble qu'on ait franchi la limite et qu'on amorce le mouvement de descente : l'axe de symétrie (11,2305) se situe en deçà de la dernière valeur de *x* (11,5960). Comme en réalité l'entropie n'a cessé de monter (mais très faiblement), il s'agit plutôt d'une imperfection de l'ajustement.



Dans une fonction du second degré, ce point qui marque l'axe de symétrie est le quotient  $-b/2a$  qui dans nos données prend une valeur proche de 11. Ce point virtuel se situe plus ou moins loin (de 11,01 à 12,40) et cela suffit à déterminer si la trajectoire est plus raide ou plus molle, c'est-à-dire si le discours garde une énergie potentielle ou s'il a épuisé ses réserves. Le classement obtenu sous ce rapport est indiqué dans la quatrième colonne inférieure du tableau 7. Il est d'ailleurs très semblable à celui que permettent les paramètres  $a$  et  $b$ , en se rangeant toujours du côté de  $a$  lorsqu'il y a quelque désaccord entre  $a$  et  $b$  (et même en ajoutant une surenchère à la tendance marquée par  $a$ ). Étienne Évrard avait donc raison de privilégier le paramètre  $a$  et de l'interpréter comme le signe de l'essoufflement lexical. Mais l'axe de symétrie, qui tient compte à la fois de  $a$  et de  $b$ , nous paraît un meilleur indice. Et surtout, il ne nous semble pas que le paramètre  $b$  ait quelque rapport direct avec la diversité lexicale. L'indice de la diversité (ou richesse) lexicale, c'est la valeur elle-même de l'entropie<sup>12</sup>, à tel moment d'un texte que l'on voudra et particulièrement à la fin, quand les jeux sont faits. La forme de la courbe est une autre notion, qui met en action le mouvement du texte et indique une tendance. Ainsi, un homme peut être riche ou pauvre (c'est l'entropie) et, étant ce qu'il est, s'employer à l'être plus ou à l'être moins (c'est la tendance).

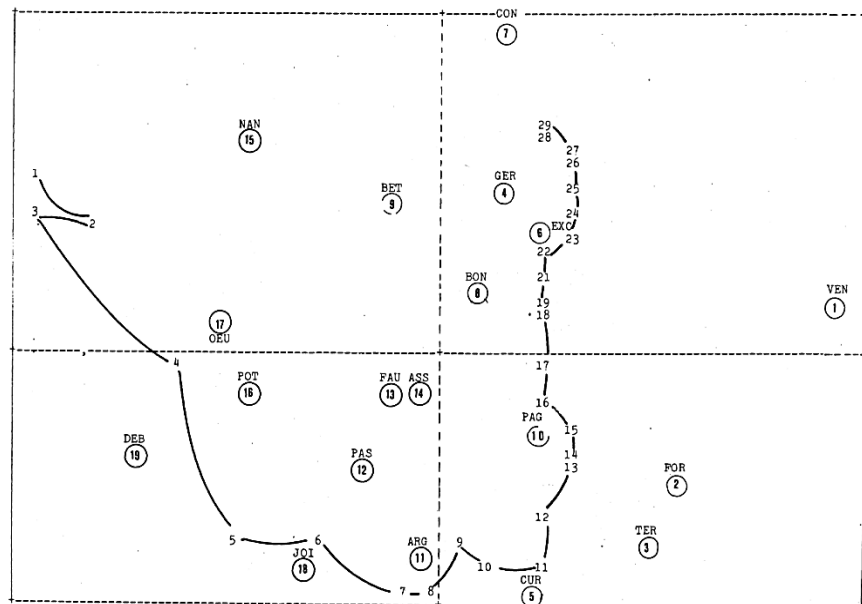
Mais le tableau à deux dimensions comme le tableau 5 ne peut pas ne pas donner l'idée d'une analyse factorielle. De si nombreux exemples de ce puissant outil ont été proposés dans le domaine linguistique qu'il semble oiseux d'en exposer les principes. Nous ne parlerons ici que des contraintes, auxquelles il a fallu souscrire en écartant les lignes et les colonnes trop courtes afin d'avoir une matrice exempte de trous. On n'a donc retenu que les 29 premières lignes, en considérant que tous les textes étaient suspendus à la 290<sup>e</sup> page, ce qui est suffisant pour l'indication d'une tendance. Les trois textes qui n'atteignaient pas cette limite (*Thérèse Raquin*, *Madeleine Férat*, et *le Rêve*) n'ont pas été pris en compte. C'est donc un tableau de 29 lignes (29 tranches de 10 pages) et de 19 colonnes (19 textes) qui a été soumis à une analyse factorielle des correspondances. Le programme reconnaît sans peine l'existence de données sérielles et sur le graphique 9, il dispose dans l'ordre les 29 maillons d'une chaîne qui décrit un croissant, comme c'est l'habitude : dans le quadrant supérieur gauche, les 3 premières tranches, puis les 5 suivantes dans le quadrant inférieur gauche, les 9 suivantes de l'autre côté

---

<sup>12</sup> Avec les réserves que nous avons dites.

de l'axe des y et les 12 dernières en haut et à droite. Or, sur un tel plan les textes se répartissent suivant que leur entropie est relativement plus forte dans les premières ou dans les dernières tranches. Ceux qui dépensent dès les premières pages leurs ressources lexicales sont à gauche, ceux qui les ménagent et les réservent pour la fin, à droite. La répartition selon le 1<sup>er</sup> facteur (ce qui sépare la gauche de la droite sur le graphique 9) reprend clairement le classement auquel nous a conduit l'axe de symétrie dans la méthode précédente. Les deux termes extrêmes sont les mêmes et sur la marge droite, le *Ventre de Paris* s'oppose à la *Débâcle* à l'extrême gauche. Nous avons fait figurer en médaillon le classement de la méthode précédente. On voit qu'il est partout respecté sauf une entorse mineure dans sa partie centrale, la *Bête Humaine* et une *Page d'Amour* (rangs 9 et 10) ayant passé la frontière.

Figure 9 : Analyse factorielle de l'entropie (19 textes et 29 tranches)



Observons ici encore qu'il s'agit de tendances et que l'analyse porte sur des profils, la valeur absolue de l'entropie n'intervenant pas. C'est pourquoi l'analyse ignore les riches et les pauvres et ne veut considérer que ceux qui s'enrichissent et ceux qui s'appauvrissent. Et de fait, on trouve des riches et des pauvres des deux côtés. Notons que l'« essoufflement » lexical – le tassement de l'entropie – se produit plus



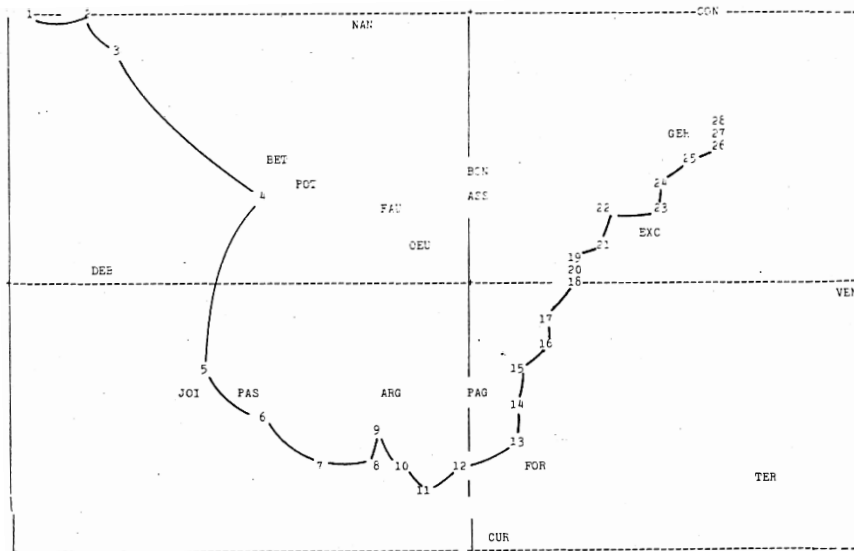


Figure 11 : Analyse factorielle du rapport  $V/\sqrt{N}$  (19 textes et 29 tranches)

Comment ne pas voir dans ce graphique la pure superposition de celle de la figure 9, qui avait été réalisée à partir de l'entropie ? Même forme en croissant ordonnant les tranches de la 1<sup>ère</sup> à la 29<sup>e</sup>. Même disposition des textes : à gauche ceux dont la variété lexicale s'épuise plus vite, à droite ceux qui renouvellent davantage leur vocabulaire.

Ainsi retrouvons-nous Pierre Guiraud au terme de notre parcours. La formule si simple qu'il a proposée jadis et qu'on a parfois décriée donne au bout du compte les mêmes indications que les calculs les plus sophistiqués – que Guiraud avait aussi suggérés. N'a-t-on donc pas fait de progrès depuis vingt ans dans une discipline que Guiraud a fondée pour l'abandonner une fois qu'elle fut devenue grande ? Certes on a gagné en puissance et les dénombrements qu'il a faits, lui, à la main, se font maintenant à grande échelle et à grande vitesse avec des machines. On a gagné aussi en précision et en fiabilité avec la facilité des contrôles et des regroupements. Mais n'a-t-on pas perdu en puissance de raisonnement, en fertilité d'imagination, en rapidité d'anticipation ? La longueur et l'inconfort des calculs manuels ont précipité Guiraud sur la voie des raccourcis, des audaces, des hypothèses et des découvertes.

Mais hélas, quand sont venus les ordinateurs, Guiraud s'en est allé.