

L'ANALYSE DE LA FREQUENCE DES ENCHAINEMENTS

Munie des arguments présentés dans le tableau ci-dessous, la macro-expression « chaines » recherche, dans le tableau d'entrée *filin*, les chaînes de deux mots consécutifs puis calcule le nombre, la fréquence, le rang et l'indice Pareto (produit du rang par la fréquence) de chaque chaîne de caractères.

Argument	Signification
ligne	Nom symbolique Sas d'une variable chaîne de caractères.
filin	Nom symbolique du tableau d'entrée (Data) Sas qui contient la chaîne de caractères <i>ligne</i> .
filout	Nom symbolique du tableau (Data) Sas de sortie.

Dès le début de son exécution, la macro-expression « chaines » transforme au moyen de la fonction TRANSLATE ci-dessous toutes les lettres accentuées en lettres non accentuées majuscules. La forme conjuguée « a » de l'auxiliaire « avoir » est ainsi assimilée à la préposition « à ». Il est évident que vous pouvez en décider autrement en modifiant l'ordre TRANSLATE.

```
%MACRO chaines (ligne,filin,filout);
DATA travail;
SET &filin;
%traduis(&ligne);
&ligne=TRANSLATE(&ligne , 'AAAEIEEEIIIOUUU', 'äääéèëîïöùü');
&ligne=UPCASE(&ligne);
i=1;
debut : mot = SCAN (&ligne,i,"$()&:~!.,;/? ");
mot = COMPRESS(mot);
mot=TRANSLATE (mot,'000000000', '123456789');
IF INDEX(mot,'0') NE 0 THEN GOTO finligne;
IF mot=' ' THEN GOTO finligne;
ELSE DO;
  molag=LAG1(mot);
  chaine = compress(molag !! '_' !! mot) ;
  OUTPUT;
  i=i+1;
  GOTO debut ;
END;
finligne :
PROC FREQ NOPRINT ;
TABLES chaine /OUT=travail ;

/* Ajout du rang et de l'indice pareto ; */

PROC SORT DATA=travail ;
BY DESCENDING COUNT;
DATA &filout ;
SET travail ;
RETAIN numrang 0;
IF COUNT NE LAG1(COUNT) THEN numrang=numrang+1;
pareto=COUNT*numrang;
%MEND ;
```

Exemple d'utilisation de la macro-expression « chaines »

Dans le tableau d'entrée Sas *madata* qui contient un flux discursif sur l'actualité internationale de janvier 2001 à octobre 2001, il existe une variable chaîne de caractères « textes » à laquelle on peut appliquer la macro-commande %chaines comme suit :

```
%chaines (textes, madata, resultat) ;
RUN ;
```

A l'image du tableau (Data) de sortie *resultat ci-dessous*, les chaînes de caractères rencontrées sont présentées par nombre décroissant (COUNT) avec la fréquence POURCENT, le rang NUMRANG et l'indice PARETO de chaque chaîne de caractères. Toutes les chaînes n'ont pas la même signification linguistique. Ainsi, les chaînes comme DE LA et A LA relèvent de la nature de la langue française alors qu'une chaîne comme « Jacques Chirac » relève du langage. Une éventuelle transformation du langage ne peut être détectée après une seule exécution de la macro-commande « chaines » : deux exécutions sont au minimum nécessaires pour détecter l'éventuelle évolution du langage.

The SAS System					
					1
11:45 Friday, February 1, 2002					
OBS	CHAINE	COUNT	PERCENT	NUMRANG	PARETO
1	DE LA	356	0.96642	1	356
2	DE L	336	0.91213	2	672
3	A LA	168	0.45606	3	504
4	AUJOURD HUI	137	0.37191	4	548
5	ET PUIS	127	0.34476	5	635
6	A L	102	0.27690	6	612
7	BONJOUR A	96	0.26061	7	672
8	L ACTUALITE	89	0.24160	8	712
9	C EST	88	0.23889	9	792
10	A TOUS	83	0.22532	10	830
11	ETATS UNIS	80	0.21717	11	880
12	EN FRANCE	78	0.21174	12	936
13	CE MATIN	77	0.20903	13	1001
14	DANS LE	74	0.20088	14	1036
15	IL Y	72	0.19546	15	1080
16	Y A	72	0.19546	15	1080
17	LE PRESIDENT	69	0.18731	16	1104
18	A ETE	68	0.18460	17	1156
19	DANS L	64	0.17374	18	1152
20	DANS LA	63	0.17102	19	1197
21	D UN	60	0.16288	20	1200
22	SUR LE	58	0.15745	21	1218
23	DANS UN	56	0.15202	22	1232
24	DANS LES	55	0.14931	23	1265
25	PROCHE ORIENT	54	0.14659	24	1296
26	LA UNE	52	0.14116	25	1300
27	SUR LA	49	0.13302	26	1274
28	UNE DE	48	0.13030	27	1296
29	AU PROCHE	47	0.12759	28	1316
30	HIER SOIR	45	0.12216	29	1305
31	JACQUES CHIRAC	44	0.11945	30	1320
32	SUR L	44	0.11945	30	1320
33	S EST	44	0.11945	30	1320
34	CE JOURNAL	43	0.11673	31	1333
35	POUR LA	42	0.11402	32	1344
36	PREMIER MINISTRE	42	0.11402	32	1344
37	SUR LES	41	0.11130	33	1353
38	D UNE	40	0.10859	34	1360
39	LA FRANCE	40	0.10859	34	1360
40	DANS UNE	38	0.10316	35	1330

L'ordre trouvé peut ensuite être conforté ou infirmé en étudiant la liaison entre le nombre et le rang de chaque chaîne de caractères avec la procédure Proc REG; Model Count=Numrang ;