

RECUENTO ESTADISTICO DE PALABRAS, LETRAS, DIGRAMAS Y TRIGRAMAS EN TITULOS DE ARTICULOS CIENTIFICOS Y TECNICOS EN ESPAÑOL

Francisco Gutiérrez Muñoz*, Gloria del Rey Gutiérrez; Alfredo del Rey Guerrero*

Resumen: Se muestran los resultados de una serie de estudios estadísticos realizados con el fin de determinar las frecuencias de aparición de las palabras, letras, digramas y trigramas contenidos en una colección de 12.540 títulos de artículos científicos y técnicos en español obtenidos a partir de la base de datos ICYT. Se discute la metodología seguida, y se sugieren posibles aplicaciones de tales resultados.

Palabras clave: Digramas, Trigramas, títulos, lingüística estadística.

Abstracts: A series of statistic studies were carried out to determine the occurrence frequency of words, characters, digrams and trigrams, in a total of 12.540 titles of scientific and technic papers indexed in the database ICYT. The are discussed the used methods and are exposed the probable uses of the obtained results.

Keyword: Digrams, trigrams, titles, statistics linguistic.

Introducción

La base de datos del Índice Español de Ciencia y Tecnología desarrollada en el ICYT plantea problemas de ocupación excesiva de memoria en disco, lo que hace conveniente su compactación. También, los errores ortográficos introducidos durante la grabación de los registros dan lugar a una pérdida de calidad en el producto resultante y a un aumento notable del tamaño de los ficheros de descriptores que se han de generar para la recuperación automatizada de información a partir de dicha base de datos. En consecuencia, se hace necesario dar solución a tales problemas.

Con tal fin, decidimos que lo más conveniente sería desarrollar, nosotros mismos, los programas de compresión de textos y de detección y corrección de errores ortográficos y tipográficos que se precisan. De entre los métodos existentes orientados a tal fin (1-6), decidimos adoptar, por su sencillez y eficacia, los que se basan en estudios de frecuencias de aparición de digramas y trigramas (combinaciones de dos o tres caracteres, respectivamente).

Por ello, y con el objeto de disponer de una base sólida para poder desarrollar tales programas, utilizando el miniordenador SECOINSA 40 del ICYT y el lenguaje COBOL, llevamos a cabo una serie de recuentos estadísticos a partir de una colección de 12.540 títulos de artículos científico-técnicos escritos en español, que se obtuvieron de la base de datos del Índice Español de Ciencia y Tecnología. Parte de los resultados obtenidos ya fueron publicados (7).

En este trabajo se reflejan, de forma resumida, los resultados más relevantes de todos los recuentos estadísticos que se han realizado en la fase previa al desarrollo de los programas de compresión de textos y de corrección de errores, actualmente en marcha.

* Instituto de Información y Documentación en Ciencia y Tecnología (ICYT), Madrid. CSIC.

Resultados

La colección de partida contiene 12.540 títulos de artículos científicos y técnicos originalmente en español. Dado que dichos títulos habían sido grabados utilizando solamente letras mayúsculas, no se han tenido en cuenta, en este estudio, ni minúsculas, ni acentos ni signos diacríticos.

La distribución temática de los títulos incluidos en la colección de trabajo se muestra en la Tabla 1, y es un reflejo de la cobertura del propio Índice Español de Ciencia y Tecnología, que abarca todos los campos de la ciencia y la tecnología con excepción de la medicina, cubierta, de manera exhaustiva, por otra publicación similar: el Índice Médico Español.

Tabla 1

Distribución temática de títulos

Area Temática	Incidencia (%)
Tecnología	37.66
Ciencias Agrarias	15.12
Ciencias Biológicas	14.72
Ciencias Químicas	9.5
Ciencias de la Tierra	9.38
Ciencias Matemáticas	5.52
Ciencias Físicas	5.32
Farmacia	2.76

A partir de dicha colección de títulos se independizaron en total 142.944 palabras, de las cuales sólo 18.099 eran diferentes. El promedio de palabras extraídas por título fue de 11'3, con valores extremos comprendidos entre 1 y 48 palabras por título. Se consideró como palabra a cualquier secuencia de caracteres alfanuméricos que comenzase por una letra y que estuviese delimitada por espacios en blanco, signos de puntuación o caracteres especiales. Su frecuencia de aparición resultó ser muy variable; así: un 54,44 por ciento de las palabras diferentes aparece una sola vez en la colección; por el contrario, un 0,56 por ciento de tales palabras, 101 palabras, aparece con mucha frecuencia, suponiendo su aportación el 50,04 por ciento del total de ocurrencias. En la Tabla 2 pueden verse las 20 palabras más frecuentes: su aportación supone el 40,17 por ciento del total; son consideradas como palabras vacías.

La longitud media de las palabras contenidas en los títulos resultó ser de 5,81 caracteres. Sin embargo, si sólo se tiene en cuenta la colección de palabras indentificadas como diferentes, la longitud media resultante es de 8,97 caracteres por palabra. En cualquier caso, la gran mayoría de las palabras tiene una longitud menor de 16 caracteres, el 99,28 por ciento de las palabras de los títulos, frente al 96,86 de las palabras diferentes.

Partiendo de la colección de palabras diferentes obtenida a partir de los títulos y teniendo en cuenta sus frecuencias de aparición respectivas, se es-

tudiaron las distribuciones de frecuencias de letras, digramas y trigramas. Todas ellas reflejan la desigual frecuencia de aparición de las palabras en los títulos de la colección de partida y, en particular, la fuerte incidencia que tienen las palabras de alta frecuencia incluidas en la Tabla 2.

Tabla 2

Relación de palabras de alta frecuencia de aparición		
Palabra	Frecuencia	Longitud
DE	20.227	2
LA	6.781	2
EN	5.693	2
Y	4.407	1
DEL	3.353	3
EL	2.362	2
LOS	2.304	3
LAS	1.755	3
SOBRE	1.314	5
PARA	1.263	4
ESTUDIO	1.260	7
A	1.142	1
CON	1.015	3
POR	944	3
UN	821	2
I	611	1
AL	607	2
UNA	572	3
II	563	2
SU	431	2

Las frecuencias de aparición de las letras simples en los títulos de la colección se muestran en la Tabla 3. Puede apreciarse en ella la preponderancia de las vocales, cuya presencia supone el 45,23 por ciento de los caracteres más frecuentes, destacan las letras; S, R, N, C, L, D y T.

Para determinar la ocurrencia de digramas y trigramas se consideró que a cada palabra de la colección de palabras diferentes le precede un espacio en blanco y le sigue otro blanco. Los digramas y trigramas podrían estar formados por combinaciones de caracteres que estuviesen contenidos en una lista de 38, integrada por las 27 letras simples del alfabeto español, los diez números dígitos y el espacio en blanco. En las tablas de digramas y trigramas (tabla 4 y tabla 5, respectivamente) el espacio en blanco se ha representado mediante un guión. Por razones de simplificación, en esta primera aproximación al tema no se han tenido en cuenta como posibles integrantes de digramas y trigramas a los signos ortográficos más frecuentes, como son: el punto, la coma, el punto y coma, etc., conscientes de que aunque ello no era interesante desde el punto de vista de la corrección de errores sí lo era desde el de la comprensión de textos.

Tabla 3

Frecuencia de aparición de letras

Letra	Frecuencia	Letra	Frecuencia
E	99.789	G	9.401
A	95.304	B	9.305
O	79.085	F	8.759
I	77.641	V	7.551
S	60.264	Y	5.565
N	56.622	H	4.087
R	50.933	Z	30.033
C	50.696	X	2.146
L	49.133	Q	1.754
D	48.592	J	1.585
T	40.944	Ñ	1.351
U	23.934	K	352
M	21.979	W	216
P	19.974		

Para hacer el recuento de los digramas se creó una tabla de 38 filas por 38 columnas, en ellas, cada uno de sus elementos correspondía al contador de frecuencias de un determinado digrama. Por limitaciones del ordenador con el que se ha trabajado, en el caso de los trigramas hubo que seguir un procedimiento diferente, que consistía en simular el desplazamiento, sobre cada línea de texto, de una ventana de tres celdillas (una por cada letra), con lo que se pudieron independizar, para posterior recuento, todos los trigramas contenidos en la colección de palabras.

Aunque en teoría podrían existir 1.444 digramas diferentes (38 x 38), sólo se encontraron 769, de los cuales 33 aportan el 50 por ciento de los 972.789 digramas obtenidos. Los 200 que se muestran en la Tabla 4 aportan el 96,12 por ciento del total de digramas. En dicha tabla, de cada digrama se indica su frecuencia de aparición y su versatilidad (entendida ésta como el número de palabras diferentes en las que aparece).

El número de trigramas posibles sería en nuestro caso de 54.872 (38 x 38 x 38). Sin embargo, se han encontrado solamente 5.628 trigramas diferentes, de los cuales 200 aportan el 51,87 por ciento del total de los 756.646 trigramas independizados. En cambio, hay otros 4.000 trigramas cuya aportación solamente supone el 4,32 por ciento de todos los trigramas independizados. La Tabla 5 muestra la frecuencia de aparición y versatilidad de los 200 trigramas más frecuentes.

Discusión y Conclusiones

En este estudio se ha partido de una colección de títulos grabados utilizando sólo letras mayúsculas. De haber sido escritos empleando mayúsculas y minúsculas, ¿en qué medida podrían haberse visto afectados los resultados

finales que se han mostrado? En relación con la ocurrencia global de letras, en nada, si bien, con respecto a las mayúsculas, podría comprobarse la incidencia destacada de algunas de ellas, por ser iniciales de palabras que frecuentemente aparecen al comienzo de un título o de nombres propios de uso frecuente.

En relación con digramas y trigramas, supondría la posibilidad de encontrar una mayor variedad de ellos, al poder intervenir mayúsculas y minúsculas en su composición; sin embargo, aunque las ocurrencias de parte de los digramas considerados en nuestro estudio se verían afectadas, no es probable que se introdujeran alteraciones importantes en las listas de digramas y trigramas más frecuentes (Tablas 4 y 5), si bien se podría apreciar en ellas la presencia de digramas y trigramas que comiencen con mayúsculas.

De otra parte, al tener que trabajar con títulos en mayúsculas, tampoco se ha podido obtener datos sobre la repercusión que pudieran tener los acentos y signos diacríticos que, en el caso de los acentos, si creemos que pudiera ser significativa, ya que van ligados a letras de muy alta frecuencia, como son las vocales, y a gran número de palabras diferentes, muchas de las cuales son de frecuente aparición.

Los resultados obtenidos tienen un valor orientativo acerca de lo que pueda encontrarse al considerar todos los títulos incluidos en los 50.000 registros que en la actualidad contiene la Base de Datos de Información en Ciencia y Tecnología (ICYT). No pretenden ser extrapolables a otros tipos de textos; sin embargo, parte de dichos resultados si pueden ser útiles para quienes trabajen en el desarrollo o explotación de bases de datos bibliográficas de cualquier tipo, y para quienes hayan de desarrollar índices permutados basados en títulos.

Tabla 4

D	S-	32.917	6.391	5.793	1.054	VA	2.293	462	213
F	-D	28.049	972	5.775	1.400	-C	2.228	576	387
D	DE	27.854	940	5.680	1.494	EP	2.227	704	243
F	E-	25.966	989	5.464	1.522	-U	2.151	381	243
D	A-	25.602	3.648	5.435	1.209	SU	2.184	172	471
F	N-	18.081	1.296	5.338	1.016	FI	2.181	358	390
D	ON	17.544	2.549	5.206	1.044	CU	2.152	532	439
F	-E	16.727	1.210	5.093	1.063	LU	2.162	415	380
D	EH	16.553	2.240	4.992	1.016	UL	2.087	626	375
F	ES	16.042	2.662	4.943	1.248	-R	2.059	716	213
D	O-	15.936	2.863	4.801	1.353	OD	2.032	396	315
F	CI	15.833	2.034	4.685	1.353	OG	2.017	435	65
D	IO	15.708	2.038	4.591	1.353	-O	1.988	448	481
F	OS	14.584	2.680	4.588	1.177	OG	1.968	330	365
D	LA	13.525	1.532	4.508	1.177	PL	1.965	446	343
F	-L	13.203	1.632	4.495	1.177	VI	1.959	432	254
D	-C	12.544	2.185	4.322	1.177	IV	1.942	479	139
F	AS	12.213	2.495	4.318	1.177	RM	1.937	388	257
D	RA	11.372	2.222	4.240	1.064	I-	1.911	346	313
F	IC	10.915	2.775	4.222	1.064	US	1.911	567	394
D	CO	10.873	1.818	4.219	1.064	-H	1.896	668	403
F	-A	10.820	1.618	4.182	1.064	HU	1.891	363	315
D	AC	10.493	1.680	4.178	1.064	RT	1.881	260	260
F	AL	9.961	1.897	4.178	1.064	OG	1.811	363	157
D	ER	9.899	2.247	4.164	1.064	DU	1.805	559	224
F	L-	9.814	609	4.106	1.064	MP	1.789	322	299
D	CA	9.315	2.051	4.092	1.064	GI	1.774	453	50
F	EL	8.733	1.693	4.082	1.064	QU	1.760	867	155
D	NT	8.590	1.571	4.004	1.064	OT	1.731	373	280
F	IN	8.337	2.029	3.905	1.064	D-	1.723	313	289
D	RE	8.318	1.561	3.815	1.064	OC	1.715	509	169
F	RI	7.885	2.039	3.733	1.064	GE	1.694	429	309
D	RO	7.856	1.818	3.615	1.064	CC	1.691	159	283
F	OR	7.856	1.669	3.604	1.064	CC	1.642	130	233
D	AR	7.339	1.691	3.573	1.064	UD	1.574	488	116
F	TO	7.323	1.975	3.521	1.064	BE	1.574	384	278
D	-S	7.323	1.132	3.521	1.064	RR	1.521	363	169
F	TA	7.038	1.676	3.521	1.064	VE	1.512	315	271
D	ST	6.948	1.828	3.521	1.064	GA	1.509	384	285
F	TI	6.841	1.828	3.521	1.064	AP	1.489	272	302
D	IA	6.452	1.828	3.521	1.064	GR	1.480	364	220
F	AN	6.452	1.828	3.521	1.064	XZ	1.438	381	245
D	NA	6.431	1.828	3.521	1.064	UI	1.428	338	155
F	LO	6.431	1.828	3.521	1.064	UA	1.410	223	49
D	AD	6.348	1.774	3.521	1.064	FO	1.408	420	280
F	TD	6.348	1.774	3.521	1.064	AB	1.398	399	193
D	TY	6.322	1.748	3.521	1.064	VO	1.398	253	138
F	LY	6.322	1.748	3.521	1.064	OP	1.385	506	64
D	-N	6.010	1.280	3.521	1.064				

Tabla 5

T	F	V	T	F	V	T	F	V	T	F	V
-DE	25.337	377	NES	2.287	387	STU	1.460	24	TUR	1.058	171
DE-	20.323	38	TIC	2.265	567	IDO	1.442	341	HAL	1.051	189
OB-	12.440	2.068	-NE	2.237	285	OBR	1.436	35	QUI	1.049	229
ION	11.315	1.270	ONE	2.200	354	TUD	1.428	17	INT	1.047	212
ON-	10.443	998	NCI	2.172	280	-AC	1.425	211	ATO	1.028	308
CIO	10.306	1.072	-BO	2.075	134	UDI	1.414	30	RAB	1.027	252
AS-	9.925	1.880	WFO	2.047	310	MEN	1.411	282	TAC	1.014	171
-LA	9.014	1.75	RIA	2.039	395	ARI	1.391	420	OLO	1.005	310
ES-	8.271	1.619	IEM	2.026	342	MIC	1.384	343	-RO	1.004	284
LA-	7.488	881	ALI	2.001	357	RIC	1.367	368	LO-	1.001	287
ACI	7.349	881	TOB	1.996	328	CTO	1.364	178	ERM	996	177
-EN	6.831	184	ALK	1.989	413	TCI	1.349	267	RIO	994	235
-CO	6.034	725	-PO	1.980	256	SOB	1.346	27	DA-	992	385
EN-	6.024	94	CO-	1.979	544	DES	1.337	287	OMP	981	133
EL-	5.953	45	NA-	1.967	500	AD-	1.337	231	UCI	980	69
ICA	4.951	1.082	DOB	1.939	483	NER	1.335	167	UCI	974	132
ENT	4.640	764	DAD	1.927	315	ACT	1.330	213	SIO	972	169
-ES	4.005	364	EDI	1.924	483	TES	1.327	234	MAB	959	271
IA-	3.654	654	TEI	1.916	465	USA	1.327	119	LAC	957	173
DEL	3.608	48	DIO	1.881	119	NRK	1.281	134	CID	955	161
CON	3.542	362	RAC	1.849	254	STE	1.271	228	INE	945	151
ICC	3.504	1.053	TRA	1.845	405	TAL	1.266	220	CAL	947	179
EST	3.368	411	TR-	1.807	280	ANA	1.254	336	ENR	945	156
-PR	3.077	407	COB	1.784	513	ILA	1.237	182	EDI	947	179
LOS	2.987	140	COB	1.784	513	ILA	1.237	182	GEN	932	96
-SL	2.965	117	DI	1.758	329	ONT	1.237	182	GEN	930	227
NTY	2.938	607	AGU	1.758	329	ONT	1.237	182	GEN	929	399
RA-	2.913	335	ENC	1.744	226	LIC	1.235	262	MAY	921	186
-CA	2.782	872	ENC	1.683	551	YAT	1.202	216	SFA	921	50
TRN	2.766	566	ENC	1.683	551	YAT	1.202	216	SFA	914	151
AP-	2.751	459	ENC	1.683	551	YAT	1.202	216	SFA	906	247
TAB	2.693	477	ENC	1.683	551	YAT	1.202	216	SFA	906	247
IDA	2.693	477	ENC	1.683	551	YAT	1.202	216	SFA	899	213
-DE	2.592	621	ENC	1.683	551	YAT	1.202	216	SFA	899	213
ADP	2.575	505	ENC	1.683	551	YAT	1.202	216	SFA	894	38
AP-	2.533	227	ENC	1.683	551	YAT	1.202	216	SFA	893	82
-LO	2.512	72	ENC	1.683	551	YAT	1.202	216	SFA	889	35
PAR	2.465	193	ENC	1.683	551	YAT	1.202	216	SFA	883	209
CA-	2.465	511	ENC	1.683	551	YAT	1.202	216	SFA	882	8
TO-	2.391	304	ENC	1.683	551	YAT	1.202	216	SFA	881	149
ARA	2.335	308	ENC	1.683	551	YAT	1.202	216	SFA	880	171
PRO	2.319	468	ENC	1.683	551	YAT	1.202	216	SFA	879	177
INA	2.314	668	ENC	1.683	551	YAT	1.202	216	SFA	866	222
RES	2.309	494	ENC	1.683	551	YAT	1.202	216	SFA	866	190
DO-	2.299	517	ENC	1.683	551	YAT	1.202	216	SFA	864	213
CIA	2.297	321	ENC	1.683	551	YAT	1.202	216	SFA	854	68
			ATI	1.468	266	PER	1.062	274	RED	850	67
			-BI	1.468	171	ATI	1.058	348	-RU	843	56

Bibliografía

- HELD, G.; MARSHALL, T.R. *Data Compression. Techniques and Applications Hardware and Software Considerations*. New York: John Wiley & Sons, 206 p., 1987.
- CHOROS, K.; MAJEWSKA, P.; SIEMINSKI, A. *Bibliographic Database Compression. Int. Forum Inf. Doc.*, 12 (1), 28-32, 1987.
- REGHBATO, H.K. *An overview of Data Compression Techniques. Computer*, 14 (4), 71-75, 1981.
- PEPPERSON J.L. *Computer Programs for Detecting and Correcting Spelling Errors. Communications of the ACM*, 23 (12), 676-687, 1980.
- BELONOGOV, G.G.; DUGANOVA, I.S; KUZNETSOV, B.A; SMTURMAN, Ya. P.; PARTYKO, Z.V.; POZDNYAK, M.V. *Automatic Error Detection and Correction in Scientific Texts. Nauch. Tekn. Inform., Ser. 2*, 16 (3), 85-89, 1982.
- ZAMORA, A. *Automatic Detection and Correction of Spelling Errors in a Large Data Base. J. Amer. Soc. Inform. Sci.*, 31 (1), 51-57, 1980.
- GUTIERREZ MUÑOZ, F.; REY GUERRERO, A. *Estudio estadístico de palabras y caracteres en títulos de artículos científico-técnicos en español. Rev. Esp. Doc. Cient.*, 9 (2), 121-132, 1986.