

RENDICONTI *del* SEMINARIO MATEMATICO *della* UNIVERSITÀ DI PADOVA

C. MINNAJA

L. PACCAGNELLA

Variazioni dell'entropia linguistica dell'italiano scritto e calcolo di un'entropia fonematica

Rendiconti del Seminario Matematico della Università di Padova,
tome 57 (1977), p. 247-265

http://www.numdam.org/item?id=RSMUP_1977__57__247_0

© Rendiconti del Seminario Matematico della Università di Padova, 1977, tous droits réservés.

L'accès aux archives de la revue « Rendiconti del Seminario Matematico della Università di Padova » (<http://rendiconti.math.unipd.it/>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques
<http://www.numdam.org/>

Variazioni dell'entropia linguistica dell'italiano scritto e calcolo di un'entropia fonematica

C. MINNAJA - L. PACCAGNELLA (*)

§ 1. - Introduzione.

L'idea del calcolo dell'entropia di un testo linguistico è già vecchia di almeno un quarto di secolo ([1], [2]). Il calcolo effettivo è stato effettuato per alcune lingue (vd., ad es. [3]), sulla base di tabelle di frequenza dei grafemi elaborate per alcune lingue europee. È ben nota l'importanza che tali studi rivestono nella teoria matematica delle comunicazioni; in particolare, per effettuare prove di intelligibilità di un messaggio, si è ricorsi a modelli di sorgente dotati della proprietà di Markov, e simulabili mediante un elaboratore elettronico ([4]).

È invece più recente l'idea proposta da ingegneri e linguisti di considerare la lingua come un sistema isolato, e di applicare a tale sistema le considerazioni proprie dell'entropia (vd., ad es. [5], a cui è annessa un'ampia bibliografia). La proposta è stata formulata in maniera assai dubitativa, anche perchè, nel caso preso sia pur superficialmente in considerazione in [5], parrebbe che l'evoluzione di una lingua nel tempo portasse ad una diminuzione di entropia.

Evidentemente l'entropia calcolata a livello di grafemi varia a seconda del testo che si considera, ma tale variazione è in genere

(*) Indirizzo degli AA. : Istituto di Matematica Applicata, Università, Padova.

Lavoro eseguito mentre L. Paccagnella godeva di una borsa di studio del CNR, Bando n. 201/1/57 del 28/5/75.

piuttosto limitata e decresce con la lunghezza del testo scelto. L'entropia calcolata a livello di fonemi varia fortemente a seconda della trascrizione fonematica scelta. Entrambi i tipi di calcolo sono utili, a seconda se si vogliono inquadrare in uno studio per un simulatore di una sorgente markoviana o per un sintetizzatore della voce umana ([6], [7]).

Nel presente lavoro viene dapprima calcolata l'entropia di primo ordine su due testi giornalistici italiani, a livello grafematico, e il risultato viene confrontato con quelli di [3], per una verifica sperimentale dell'ipotesi espressa in [5]. Viene poi calcolata l'entropia degli stessi testi, ma con alfabeti più ampi, come proposta per modelli di sorgente markoviana più adeguati alla lingua italiana.

Successivamente viene calcolata l'entropia di primo ordine a livello fonematico, e infine viene proposta un'interpretazione diversa di un possibile calcolo dell'entropia a livello di parole.

§ 2. - Calcolo dell'entropia sull'alfabeto base.

La classica formula di Shannon

$$H_1(X) = - \sum_{x \in X} p(x) \lg_2 p(x)$$

fornisce l'entropia calcolata su un insieme X di simboli x , sul quale è stata introdotta una funzione di probabilità p . I risultati trovati da Manfrino ([3]) per H_1 sono ricavati dall'analisi di tre testi, ciascuno di 10.000 lettere dell'alfabeto italiano. I tre testi erano rispettivamente di carattere scientifico, storico e giornalistico, scritti rispettivamente nel 1955, 1940 e 1958. L'entropia media risultante era di 3,9433 bit/lettera.

La nostra analisi si esprime in un confronto tra due testi giornalistici, uno di politica interna e uno di politica estera. Entrambi i pezzi sono di 10.000 lettere, prescindendo dal fatto che queste siano soltanto le 21 lettere dell'alfabeto italiano classico. Compaiono infatti anche cifre arabe e romane (raggruppate in un unico elemento) e lettere di altri alfabeti.

Le occorrenze e le relative entropie sono riportate nella tabella I. Il testo n° 1 è preso dall'articolo di fondo del quotidiano « La Repubblica » del 24/3/1976, pagg. 1 e 3, dal titolo « La lunga notte della DC »; esso si compone di due parti, a firma rispettivamente

TABELLA I

Lettera	Testo N° 1		Testo N° 2	
	occ./10 ⁴	H_1	occ./10 ⁴	H_1
<i>a</i>	0.1191	0.3656	0.1081	0.3470
<i>b</i>	0.0077	0.0541	0.0089	0.0606
<i>c</i>	0.0441	0.1986	0.0450	0.2013
<i>d</i>	0.0403	0.0398	0.1867	0.1851
<i>e</i>	0.1089	0.3484	0.1067	0.3445
<i>f</i>	0.0095	0.0638	0.0103	0.0680
<i>g</i>	0.0195	0.1108	0.0180	0.1043
<i>h</i>	0.0096	0.0643	0.0101	0.0670
<i>i</i>	0.1069	0.3448	0.1199	0.3669
<i>l</i>	0.0648	0.2558	0.0623	0.2495
<i>m</i>	0.0213	0.1183	0.0271	0.1411
<i>n</i>	0.0743	0.2787	0.0724	0.2742
<i>o</i>	0.0977	0.3278	0.0874	0.3073
<i>p</i>	0.0268	0.1399	0.0277	0.1433
<i>q</i>	0.0039	0.0312	0.0024	0.0209
<i>r</i>	0.0652	0.2568	0.0655	0.2576
<i>s</i>	0.0513	0.2198	0.0585	0.2396
<i>t</i>	0.0728	0.2752	0.0648	0.2558
<i>u</i>	0.0276	0.1429	0.0277	0.1433
<i>v</i>	0.0132	0.0824	0.0132	0.0824
<i>z</i>	0.0131	0.0819	0.0091	0.0617
<i>j</i>	0.0000	0.0000	0.0002	0.0025
<i>k</i>	0.0000	0.0000	0.0006	0.0064
<i>w</i>	0.0001	0.0013	0.0007	0.0073
<i>x</i>	0.0000	0.0000	0.0001	0.0013
<i>y</i>	0.0000	0.0000	0.0011	0.0108
cifre	0.0023	0.0202	0.0124	0.0785
Totali	1.0000	3.9693	1.0000	4.0282

di F. De Luca e G. Valentini. Il testo n° 2 è preso dalla rubrica « Affari esteri » del settimanale « Panorama » del 23/3/1976 ; esso si compone di due parti, una intitolata « Distensione proibita » a firma di M. Conti, l'altra (fino al completamento delle 10.000 lettere) intitolata « Se il rosso vince » a firma di S. Parone. In entrambi i testi non sono stati considerati nè titoli nè sottotitoli. In appendice,

nelle tabelle I_A e II_A sono riportate le occorrenze suddivise per singole migliaia di lettere; ad esse si riferiscono i due gruppi di grafici annessi. Questi dimostrano come siano fortemente differenziate le occorrenze in ciascun migliaio di lettere.

Si può notare che le singole entropie per i tre testi calcolate in [3] erano parecchio simili tra loro, rispettivamente: 3,9453; 3,9395; 3,9453. Invece la prima metà della tabella I fornisce già 3,9693, e la seconda metà addirittura 4,0282. Per quanto l'estensione dei testi in esame sia relativa, tali differenze potrebbero già portare un supporto alla tesi che lo stile, almeno quello giornalistico, in 18 anni si è evoluto verso una maggiore entropia. In particolare è da notare la differenziazione relativamente forte tra l'entropia dei due testi attuali entrambi giornalistici, uno di politica interna e l'altro di politica estera, superiore di circa dieci volte alle differenze registrate tra i testi di [3], che pure erano di tipo piuttosto differente tra loro.

§ 3. - Ampliamenti dell'alfabeto.

I calcoli effettuati per la tabella I tendevano a raffrontare le frequenze delle varie lettere con i testi di qualche decennio fa. Qualora si verificassero, con l'evolversi della lingua, significativi spostamenti di tali frequenze, potrebbe essere utile modificare il codice telegrafico in vista di una maggiore economia del costo di trasmissione.

Se aumentiamo il numero degli elementi dell'alfabeto X l'entropia in generale crescerà. Un ampliamento abbastanza naturale è quello di considerare anche i segni di interpunzione. Nel testo n° 1 sono stati trovati 12 diversi segni di interpunzione, che hanno portato il totale del corpus a 10.330 elementi; nel testo n° 2 ne sono stati trovati 14, per un totale di 10.393 elementi. Ciò ha portato alla tabella II, nella quale le frequenze sono state rapportate a 10.000. Si noti come sia sempre maggiore l'entropia del testo n° 2, per quanto la differenza percentuale resti pressochè costante.

Questo ampliamento dell'alfabeto comporta però alcune difficoltà. Infatti, adottando un codice binario, se l'alfabeto è composto da n elementi, bisogna usare k cifre per distinguere un elemento, dove k è legato ad n dalla relazione

$$2^{k-1} + 1 < n < 2^k$$

TABELLA II

Lettera	Testo N° 1		Testo N° 2	
	freq./10 ⁴	H_1	freq./10 ⁴	H_1
<i>a</i>	0.1153	0.3593	0.1040	0.3396
<i>b</i>	0.0075	0.0527	0.0085	0.0585
<i>c</i>	0.0427	0.1942	0.0432	0.1958
<i>d</i>	0.0390	0.1826	0.0382	0.1799
<i>e</i>	0.1054	0.3422	0.1026	0.3370
<i>f</i>	0.0092	0.0622	0.0099	0.0659
<i>g</i>	0.0189	0.1081	0.0173	0.1013
<i>h</i>	0.0093	0.0627	0.0097	0.0649
<i>i</i>	0.1061	0.3434	0.1153	0.3593
<i>l</i>	0.0627	0.2506	0.0599	0.2433
<i>m</i>	0.0206	0.1155	0.0260	0.1369
<i>n</i>	0.0719	0.2731	0.0696	0.2676
<i>o</i>	0.0946	0.3218	0.0840	0.3002
<i>p</i>	0.0259	0.1367	0.0266	0.1392
<i>q</i>	0.0038	0.0304	0.0023	0.0202
<i>r</i>	0.0631	0.2516	0.0630	0.2513
<i>s</i>	0.0497	0.2151	0.0562	0.2334
<i>t</i>	0.0705	0.2697	0.0623	0.2495
<i>u</i>	0.0267	0.1396	0.0266	0.1392
<i>v</i>	0.0128	0.0804	0.0127	0.0800
<i>z</i>	0.0127	0.0799	0.0087	0.0595
<i>j</i>	0.0000	0.0000	0.0001	0.0013
<i>k</i>	0.0000	0.0000	0.0005	0.0055
<i>w</i>	0.0001	0.0013	0.0006	0.0064
<i>x</i>	0.0000	0.0000	0.0001	0.0013
<i>y</i>	0.0000	0.0000	0.0010	0.0100
cifre	0.0022	0.0196	0.0119	0.0760
.	0.0070	0.0499	0.0071	0.0507
,	0.0130	0.0813	0.0152	0.0918
;	0.0001	0.0013	0.0000	0.0000
:	0.0013	0.0121	0.0016	0.0149
?	0.0003	0.0034	0.0004	0.0045
!	0.0002	0.0023	0.0000	0.0000
()	0.0008	0.0080	0.0007	0.0073
»	0.0000	0.0000	0.0001	0.0013
« , »	0.0048	0.0372	0.0051	0.0388
—	0.0002	0.0023	0.0001	0.0013
...	0.0001	0.0013	0.0003	0.0035
+	0.0002	0.0023	0.0000	0.0000
'	0.0041	0.0323	0.0051	0.0388
%	0.0000	0.0000	0.0008	0.0082
(0),	0.0000	0.0000	0.0005	0.0055
Totali		4.1264		4.1896

ed in questo caso avremmo già bisogno di 6 cifre, che si adatterebbero ad alfabeti con un numero di elementi compreso tra 33 e 64.

Ancora maggiore è tuttavia l'entropia di un testo scritto quando si considerino diverse le maiuscole dalle minuscole. Nella nostra analisi sono state considerate *maiuscole* soltanto quelle lettere maiuscole che non comparivano in inizio di frase, in quanto una maiuscola dopo un punto è chiaramente ridondante. Ciò ha portato per il testo n° 1 ad un alfabeto di 40 simboli (22 minuscole, cifre, 17 maiuscole) e per il testo n° 2 ad un alfabeto di 46 simboli (25 minuscole, cifre, 20 maiuscole). Tali numeri sono ovviamente legati al corpus; in linea generale si dovrebbero prevedere almeno 53 simboli (26 minuscole, 26 maiuscole, cifre). Si noti che nel testo n° 2 alcune lettere compaiono più frequentemente nella forma maiuscola che nella minuscola, il che mostra come la politica estera affidi una discreta parte della sua informazione ai nomi propri. Infatti le maiuscole sono 122 nel testo n° 1 e 146 nel testo n° 2. L'entropia totale, considerando maiuscole, minuscole e cifre, senza considerare gli altri segni tipografici, è 4.0136 per il primo testo e 4.1143 per il secondo.

Un ulteriore ampliamento dell'alfabeto, comprendente maiuscole, minuscole, cifre, segni di interpunzione porta ad un codice a 7 elementi, in quanto è da prevedere che l'alfabeto sia di $53 + 14 = 67$ simboli. Le tabelle III e IV portano rispettivamente le occorrenze di tutti questi elementi e le rispettive entropie, già calcolate rapportate a 10.000, per il testo n° 1 e per il testo n° 2 (*).

Potremmo considerare che questo ampliamento dell'alfabeto è il massimo effettivamente significativo, a meno che non si voglia distinguere una vocale scritta con l'accento da una senza accento, oppure una cifra dall'altra, il che è essenziale dal punto di vista dell'informazione, ma meno importante dal punto di vista dello stile linguistico. Un discorso a parte meriterebbe lo spazio come segno tipografico.

(*) Il calcolo relativo a queste tabelle è stato eseguito da G. Zilli e F. Degan sul calcolatore Olivetti 652.

TABELLA III

Simbolo	occ.	H_1	Simbolo	occ.	H_1
<i>a</i>	1183	0.3580	<i>r</i>	648	0.2509
<i>A</i>	8	0.0080	<i>R</i>	4	0.0044
<i>b</i>	67	0.0472	<i>s</i>	510	0.2146
<i>B</i>	10	0.0097	<i>S</i>	3	0.0034
<i>c</i>	427	0.1903	<i>t</i>	727	0.2698
<i>C</i>	14	0.0129	<i>T</i>	1	0.0013
<i>d</i>	389	0.1784	<i>u</i>	276	0.1398
<i>D</i>	14	0.0129	<i>U</i>	0	0.0000
<i>e</i>	1086	0.3420	<i>v</i>	131	0.0800
<i>E</i>	3	0.0034	<i>V</i>	1	0.0013
<i>f</i>	81	0.0549	<i>z</i>	111	0.0704
<i>F</i>	14	0.0129	<i>Z</i>	20	0.0175
<i>g</i>	183	0.1032	<i>w</i>	1	0.0013
<i>G</i>	12	0.0113	<i>W</i>	0	0.0000
<i>h</i>	96	0.0628	.	72	0.0499
<i>H</i>	0	0.0000	,	134	0.0813
<i>i</i>	1068	0.3388	;	1	0.0013
<i>I</i>	1	0.0013	:	13	0.0121
<i>l</i>	646	0.2504	?	3	0.0034
<i>L</i>	2	0.0024	!	2	0.0023
<i>m</i>	204	0.1119	()	8	0.0080
<i>M</i>	9	0.0089	« »	48	0.0372
<i>n</i>	741	0.2730	—	2	0.0023
<i>N</i>	2	0.0024	...	1	0.0013
<i>o</i>	977	0.3221	—	2	0.0023
<i>O</i>	0	0.0000	'	41	0.0323
<i>p</i>	264	0.1354	cifre	22	0.0196
<i>P</i>	4	0.0044			
<i>q</i>	39	0.0304			
<i>Q</i>	0	0.0000	Totali	10.330	4.1973

TABELLA IV

Simbolo	occ.	H_1	Simbolo	occ.	H_1
<i>a</i>	1079	0.3393	<i>T</i>	1	0.0013
<i>A</i>	2	0.0024	<i>u</i>	257	0.1320
<i>b</i>	81	0.0546	<i>U</i>	20	0.0174
<i>B</i>	8	0.0080	<i>v</i>	129	0.0786
<i>c</i>	443	0.1941	<i>V</i>	3	0.0034
<i>C</i>	7	0.0071	<i>z</i>	91	0.0599
<i>d</i>	390	0.1777	<i>Z</i>	0	0.0000
<i>D</i>	8	0.0080	<i>j</i>	0	0.0000
<i>e</i>	1060	0.3360	<i>J</i>	2	0.0024
<i>E</i>	7	0.0071	<i>k</i>	2	0.0024
<i>f</i>	95	0.0619	<i>K</i>	4	0.0044
<i>F</i>	8	0.0080	<i>w</i>	1	0.0013
<i>g</i>	175	0.0992	<i>W</i>	6	0.0062
<i>G</i>	5	0.0053	<i>x</i>	1	0.0013
<i>h</i>	99	0.0640	<i>X</i>	0	0.0000
<i>H</i>	2	0.0024	<i>y</i>	11	0.0105
<i>i</i>	1197	0.3591	<i>Y</i>	0	0.0000
<i>I</i>	2	0.0024	Cifre	124	0.0762
<i>l</i>	623	0.2434	.	74	0.0508
<i>L</i>	0	0.0000	,	158	0.0918
<i>m</i>	248	0.1286	;	0	0.0000
<i>M</i>	23	0.0195	:	17	0.0151
<i>n</i>	722	0.2673	?	5	0.0053
<i>N</i>	2	0.0024	!	0	0.0000
<i>o</i>	873	0.3002	()	8	0.0080
<i>O</i>	1	0.0013	»	2	0.0024
<i>p</i>	263	0.1342	«	54	0.0394
<i>P</i>	14	0.0128	—	2	0.0024
<i>q</i>	24	0.0202	...	4	0.0044
<i>Q</i>	0	0.0000	'	54	0.0394
<i>r</i>	643	0.2484	%	9	0.0088
<i>R</i>	12	0.0113	(0),	6	0.0062
<i>s</i>	576	0.2313			
<i>S</i>	9	0.0088			
<i>t</i>	647	0.2494	Totali	10.393	4.2870

§ 4. - Entropia della lingua parlata.

Le percentuali dei fonemi calcolate in [8] ci permettono, sempre con l'uso della formula di Shannon

$$H(X) = - \sum_{x \in X} p(x) \lg_2 p(x),$$

di calcolare anche l'entropia dei singoli fonemi e quindi della lingua parlata. Usiamo qui lo stesso testo analizzato in [8] (la « Veglia d'armi » di Diego Fabbri), e la stessa trascrizione fonetica usata in [8]; tale trascrizione è molto precisa, nonostante non sia troppo convincente la motivazione per cui i fonemi z e z' vengono esclusi dall'opposizione « geminata : semplice ». Il fatto che lo spettro acustico presenti anche per z e z' semplici l'aspetto di una consonante lunga non autorizza, a nostro giudizio, a dichiarare trascurabile tale opposizione. Sono invece molto convincenti i motivi per i quali è stato scelto tale testo come esempio di lingua parlata : si tratta di un lavoro teatrale, quindi colloquiale, con argomenti di varia natura, per quanto sempre riferibile ad un linguaggio di un ceto medio-colto. Il testo conta 18.970 parole che sono state trascritte in 83.098 fonemi. La tabella V riporta i valori di H da noi ricavati e, per comodità, anche le frequenze rapportate a 10.000 come da [8]. I fonemi distinti sono 50, il che, dal punto di vista tecnico, ci porterebbe ad un codice a 6 cifre.

Come si vede, l'entropia calcolata sui fonemi è notevolmente superiore a quella calcolata sui grafemi, tuttavia non così superiore quanto potrebbe apparire considerando solo i grafemi registrati da [3].

§ 5. - Successivi livelli di entropia.

Per uno studio più avanzato si può calcolare l'entropia condizionata di livello r così definita

$$\begin{aligned} H_r(Y) &= - \sum_{x, y \in Y} p b_x(y) \lg_2 = \\ &= - \sum_{x, y \in Y} p(b_x, y) \lg_2 p(b_x, y) + \sum_{x \in X} p(b_x) \lg_2 p(b_x) \end{aligned}$$

TABELLA V

Fonema	freq./10 ⁴	H ₁	Fonema	freq./10 ⁴	H ₁
<i>a</i>	0.06432	0.25461	<i>m/</i>	0.00126	0.01213
<i>á</i>	0.03966	0.18466	<i>n</i>	0.07271	0.27496
<i>b</i>	0.00523	0.03963	<i>n/</i>	0.00122	0.01180
<i>b/</i>	0.00254	0.02189	<i>n'</i>	0.00184	0.01671
<i>č</i>	0.00744	0.05428	<i>o</i>	0.08000	0.29150
<i>č/</i>	0.00087	0.00884	<i>ó</i>	0.02271	0.12400
<i>k</i>	0.04109	0.18922	<i>ò</i>	0.01385	0.08550
<i>k/</i>	0.00205	0.01830	<i>p</i>	0.02981	0.15107
<i>d</i>	0.03313	0.15386	<i>p/</i>	0.00162	0.01501
<i>d/</i>	0.00016	0.00201	<i>r</i>	0.06837	0.26462
<i>e</i>	0.08214	0.29617	<i>r/</i>	0.00172	0.01582
<i>é</i>	0.02839	0.14588	<i>s</i>	0.04423	0.19898
<i>è</i>	0.02235	0.12255	<i>s/</i>	0.00627	0.04588
<i>f</i>	0.00825	0.05710	<i>š</i>	0.00209	0.01860
<i>f/</i>	0.00059	0.00632	<i>š</i>	0.00394	0.03147
<i>g</i>	0.00385	0.03088	<i>t</i>	0.05677	0.23495
<i>g/</i>	0.00004	0.00058	<i>t/</i>	0.00678	0.04884
<i>ǰ</i>	0.00385	0.03088	<i>ú</i>	0.01279	0.08043
<i>ǰ/</i>	0.00102	0.01013	<i>u</i>	0.00856	0.05879
<i>i</i>	0.06504	0.25642	<i>w</i>	0.01189	0.07602
<i>í</i>	0.01938	0.11025	<i>v</i>	0.02134	0.11844
<i>y</i>	0.02096	0.11687	<i>v/</i>	0.00051	0.00557
<i>l</i>	0.03196	0.15876	<i>z</i>	0.00483	0.37160
<i>l/</i>	0.00656	0.04757	<i>z'</i>	0.00020	0.00245
<i>l'</i>	0.00218	0.01927			
<i>m</i>	0.03111	0.15575	Entropia totale		5.04781

dove $p_{b_x}(y)$ è la probabilità di y condizionata dal gruppo b_x , formato da $r-1$ lettere, che precede y .

Le H_2 e H_3 per i grafemi della lingua italiana sono già state calcolate in [3]. Per i fonemi nulla finora è stato calcolato, per quanto per il calcolo di H_2 potrebbe servire di spunto la tabella II di [4], nonostante che la trascrizione sia meno fine di quella usata in [8], pur riferendosi allo stesso testo.

È interessante il calcolo dell'entropia a livello di parole. Per questo serve una valutazione piuttosto precisa della lunghezza media delle parole della lingua in questione. A livello grafematico la lunghezza media riportata in [3] è di 5,20 lettere/parola; nei testi n° 1

e n° 2 oggetto del presente studio le parole erano sempre 1911, il che conduce ad una lunghezza media di 5,23 lettere/parola. A livello di fonemi il calcolo sui dati di [8] ci conduce a 4,380 fonemi/parola.

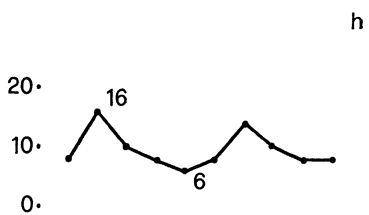
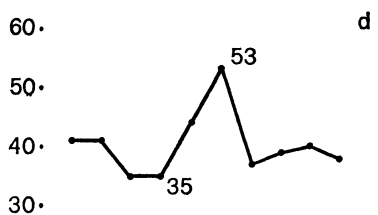
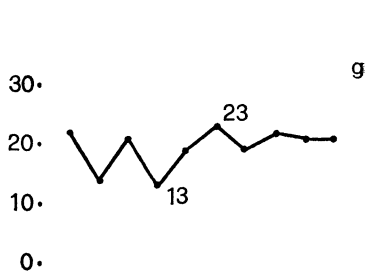
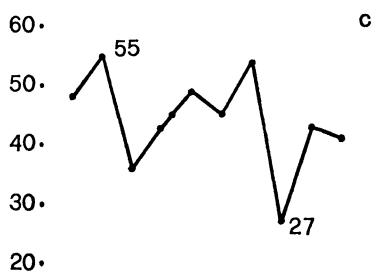
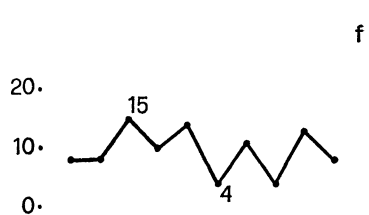
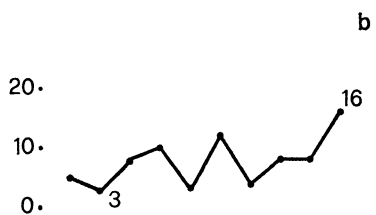
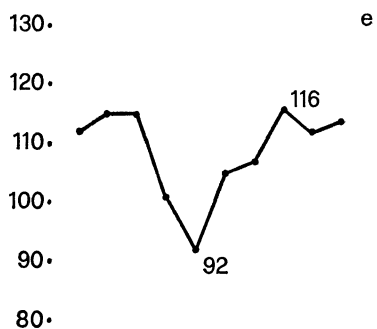
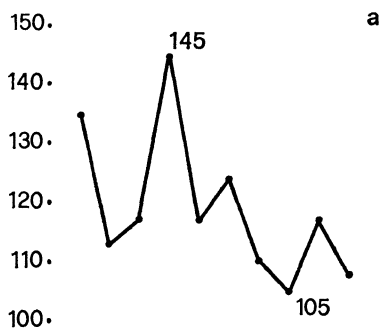
Resta aperto il problema del calcolo dell'entropia a livello di parole per un numero significativo di testi di varie epoche, per poter verificare una variazione di entropia con il tempo. Si potrebbe anche definire un vocabolario di concetti base da assumere come alfabeto: nella lingua italiana di uso comune non sono più di 700-800, nonostante che i lemmi siano molte migliaia. In base a tale alfabeto si potrebbero calcolare entropie di testi scritti e parlati, rilevando così con maggior rigore somiglianze e differenze di stile e di contenuto.

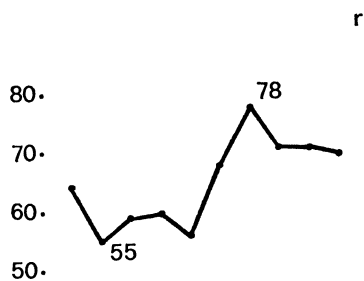
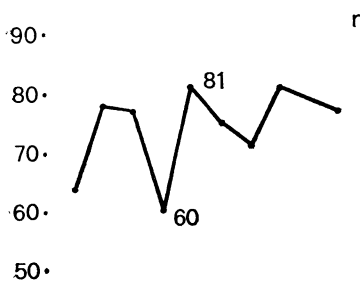
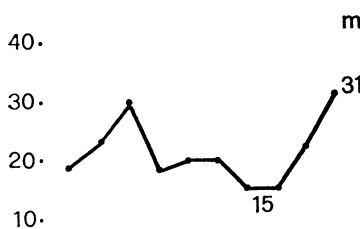
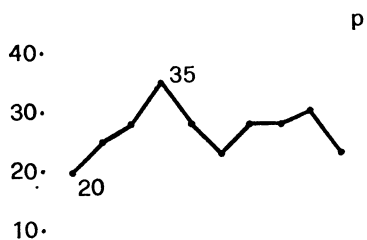
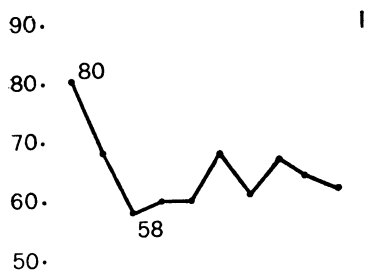
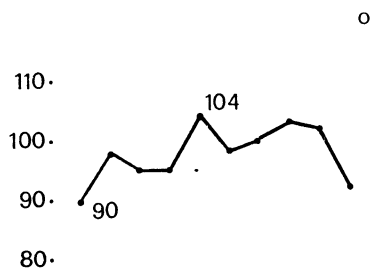
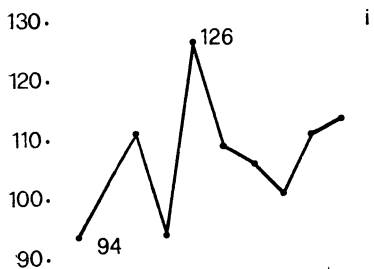
TABELLA I_A

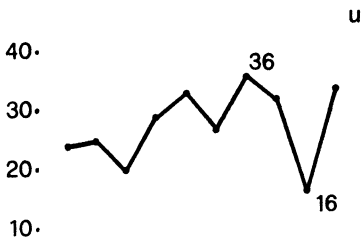
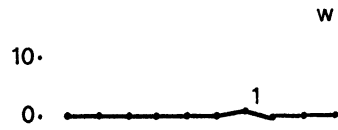
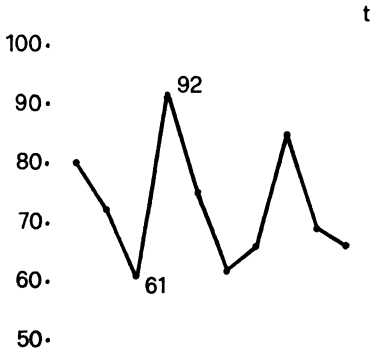
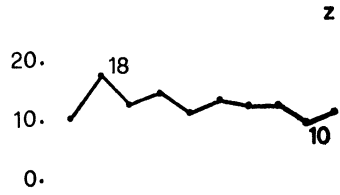
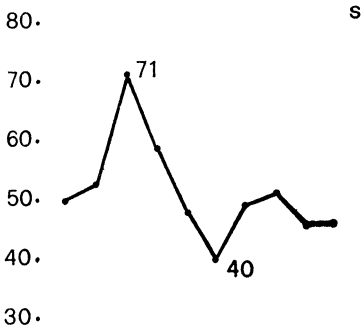
<i>a</i>	135	113	117	145	117	124	110	105	117	108
<i>b</i>	5	3	8	10	3	12	4	8	8	16
<i>c</i>	48	55	36	43	49	45	54	27	43	41
<i>d</i>	41	41	35	35	44	53	37	39	40	38
<i>e</i>	112	115	115	101	92	105	107	116	112	114
<i>f</i>	8	8	15	10	14	4	11	4	13	8
<i>g</i>	22	14	21	13	19	23	19	22	21	21
<i>h</i>	8	16	10	8	6	8	14	10	8	8
<i>i</i>	94	103	111	94	126	109	106	101	111	114
<i>l</i>	80	68	58	60	60	68	61	67	64	62
<i>m</i>	19	23	30	18	20	20	15	15	22	31
<i>n</i>	64	78	77	60	81	75	71	81	79	77
<i>o</i>	90	98	95	95	104	98	100	103	102	92
<i>p</i>	20	25	28	35	28	23	28	28	30	23
<i>q</i>	3	4	3	4	4	4	3	5	3	6
<i>r</i>	64	55	59	60	56	68	78	71	71	70
<i>s</i>	50	53	71	59	48	40	49	51	46	46
<i>t</i>	80	72	61	92	75	62	66	85	69	66
<i>u</i>	24	25	20	29	33	27	36	32	16	34
<i>v</i>	16	11	15	14	9	16	10	14	14	13
<i>z</i>	11	18	13	15	12	14	13	13	10	12
<i>j</i>	—	—	—	—	—	—	—	—	—	—
<i>k</i>	—	—	—	—	—	—	—	—	—	—
<i>w</i>	—	—	—	—	—	—	1	—	—	—
<i>x</i>	—	—	—	—	—	—	—	—	—	—
<i>y</i>	—	—	—	—	—	—	—	—	—	—

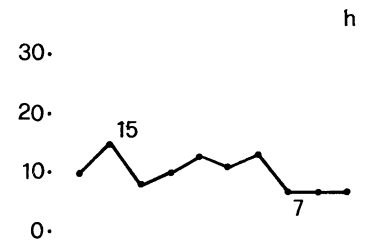
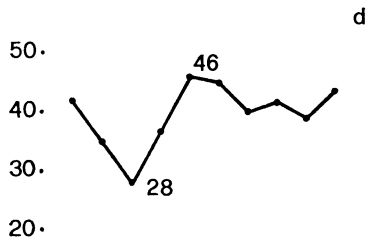
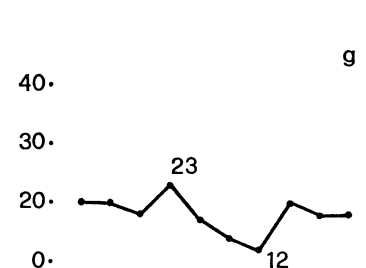
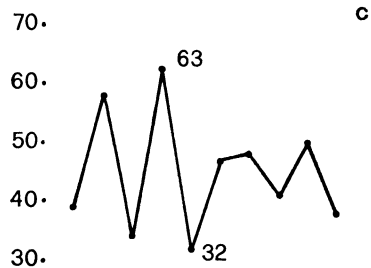
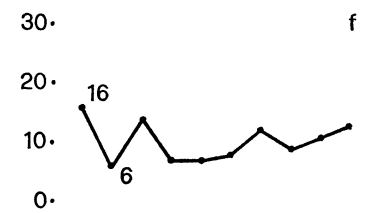
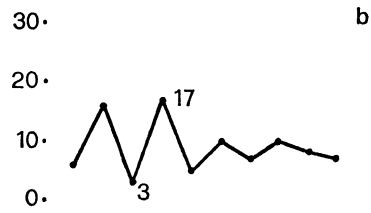
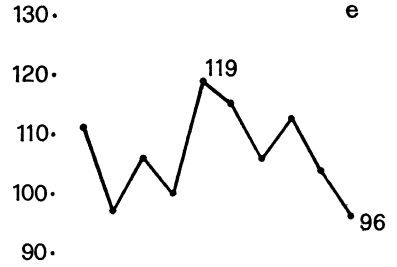
TABELLA II_A

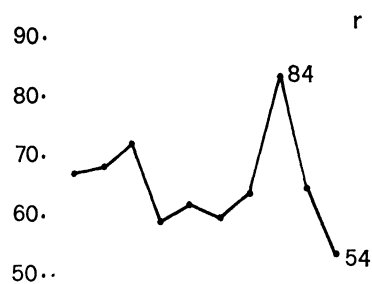
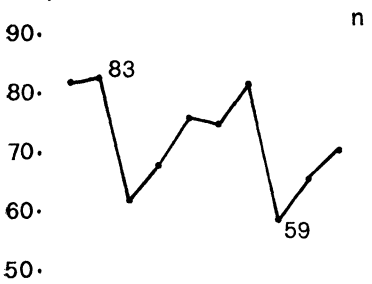
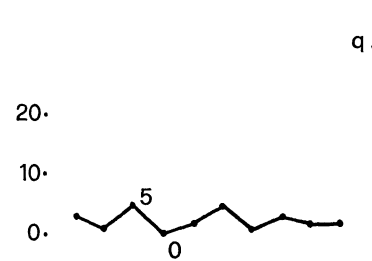
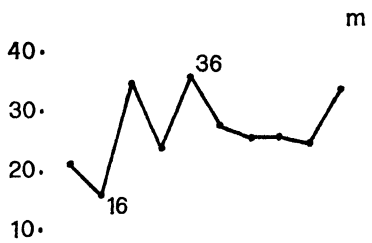
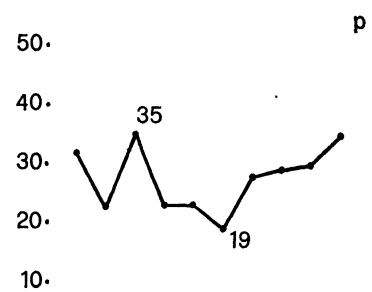
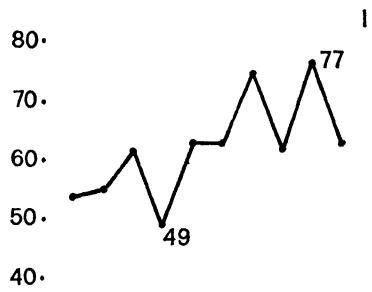
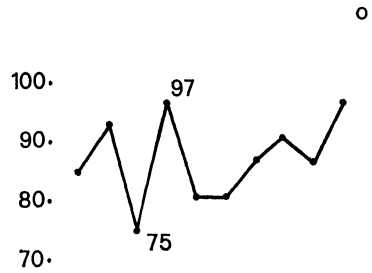
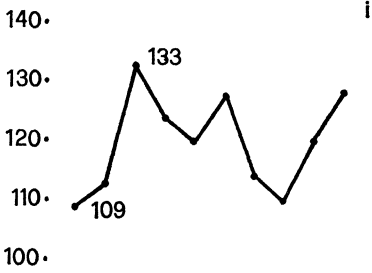
<i>a</i>	124	127	121	103	100	102	98	112	101	112
<i>b</i>	6	16	3	17	5	10	7	10	8	10
<i>c</i>	39	58	34	63	32	47	48	41	50	41
<i>d</i>	42	35	28	37	46	45	40	42	39	42
<i>e</i>	111	97	106	100	119	115	106	113	104	113
<i>f</i>	16	6	14	7	7	8	12	9	11	9
<i>g</i>	20	20	18	23	17	14	12	20	18	20
<i>h</i>	10	15	8	10	13	11	13	7	7	7
<i>i</i>	109	113	133	124	120	128	114	110	120	110
<i>l</i>	54	55	62	49	63	63	75	62	77	62
<i>m</i>	21	16	35	24	36	28	26	26	25	26
<i>n</i>	82	83	62	68	76	75	82	59	66	59
<i>o</i>	85	93	75	97	81	81	87	91	87	91
<i>p</i>	32	23	35	23	23	19	28	29	30	29
<i>q</i>	3	1	5	—	2	5	1	3	2	3
<i>r</i>	67	68	72	59	62	60	64	84	65	84
<i>s</i>	63	56	59	61	65	64	54	52	58	52
<i>t</i>	68	58	61	70	77	65	65	66	58	66
<i>u</i>	26	22	29	22	38	31	24	28	30	28
<i>v</i>	11	18	11	14	9	14	16	17	15	17
<i>z</i>	6	13	13	6	5	8	10	11	9	11
<i>j</i>	1	1	—	—	—	—	—	—	—	—
<i>k</i>	1	2	1	1	1	—	—	—	—	—
<i>w</i>	1	2	1	2	1	—	—	—	—	—
<i>x</i>	1	—	—	—	—	—	—	—	—	—
<i>y</i>	1	2	—	—	2	1	1	2	2	2

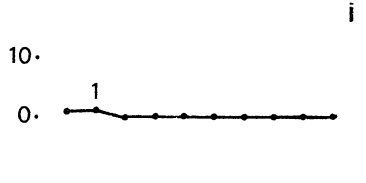
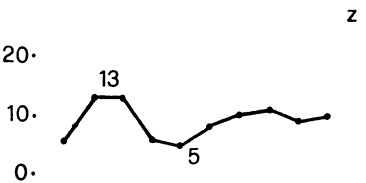
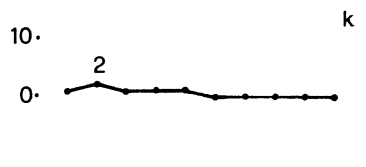
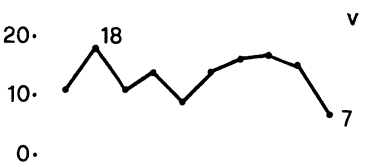
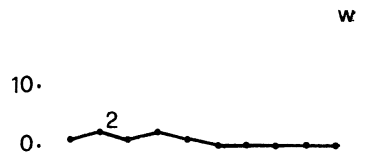
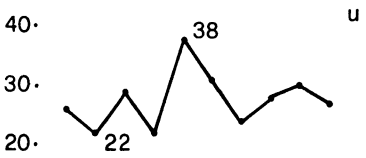
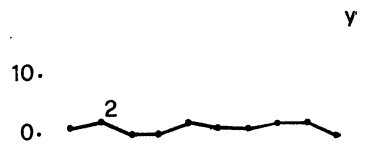
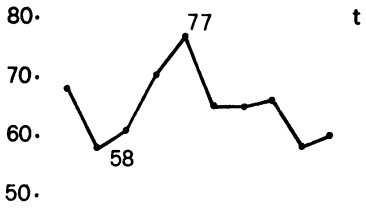
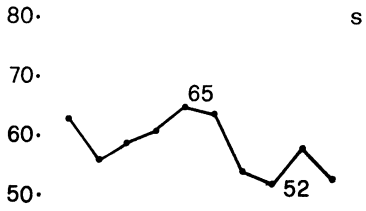












BIBLIOGRAFIA

- [1] C. E. SHANNON, *Prediction and Entropy of Printed English*, Bell S.T.J., 1951, XXX, pp. 50-64.
- [2] G. A. BARNARD, *Statistical Calculation of Word Entropies for Four Western Languages*, I.R.E. Trans. Information Theory, 1955, IT-I, p. 49.
- [3] R. MANFRINO, *L'entropia della lingua italiana ed il suo calcolo*, « Alta frequenza », 1960, XXIX, pp. 4-29.
- [4] G. B. DEBIASI - G. A. VALLI, *Sorgenti di messaggi per prove di intelligibilità della parola*, Mem. Acc. Pat. Sc. Lett. Arti, 1968, LXXX, pp. 293-314.
- [5] E. A. AFENDRAS - N. S. TZANNES - J. G. TRÉPANIÉ, *Distance, Variation and Change in Phonology: Stochastic Aspects*, Folia linguistica, 1973, 6, pp. 1-27.
- [6] G. B. DEBIASI - G. DE POLI - G. A. MIAN - C. MILDONIAN - C. OFFELLI, *Italian speech synthesis from unrestricted text for an automatic answerback system*, Proc. 8th Int. Congr. Acoustics, Londra, 1974, p. 296.
- [7] G. B. DEBIASI - G. DE POLI - G. A. MIAN - C. OFFELLI, *Voce dagli elaboratori: prospettive di sviluppo e realtà di una applicazione*, Atti 3^o Convegno di Cibernetica e Biofisica, S. Marino, 1974.
- [8] R. BUSA - C. CROATTO-MARTINOLLI - L. CROATTO - C. TAGLIAVINI - A. ZAMPOLLI, *Una ricerca statistica sulla composizione fonologica della lingua italiana parlata eseguita con un sistema IBM a schede perforate*, Proc. 12th Int. Speech & Voice Therapy Conference, Padua, 1963, pp. 542-562.

Manoscritto pervenuto in redazione il 6 aprile 1977.