

Are Blogs Edited? A Linguistic Survey of Italian Blogs

Using Search Engines

Mirko Tavosanis

Pisa University
Dipartimento di Studi Italianistici
Via del Collegio Ricci 10
I-56126 Pisa - Italy
phone: +39 050 2215065
email: tavosanis@ital.unipi.it

Abstract

Many blogs are written by people with no formal training in public writing; this could suggest a low level of editing and general correctness. A quantitative analysis of misspellings, however, shows that in their orthography Italian blogs are as well revised as conventional Italian newspaper texts. On the other hand, their editing is more careful than the editing of the average of Italian web pages.

Context: an empirical grid

The nature of the texts published on the Web is still poorly described from the linguistic viewpoint. References to the “informal” nature of all texts written for the Web can often still be found. This kind of view has been confuted in the past (see in particular Crystal 2001) but even in recent years descriptions of the linguistic features of real Web writing are scarce, even for blogs.

In this survey we will take as a starting point for an orthographic analysis an empirical grid for the description of Web texts from the linguistic point of view presented in Tavosanis (2005). The grid is intended mainly as a compendium of words to enable linguists to speak more correctly about the texts published on the Web and to place particular phenomena in context. The whole classification, slightly different from current grids, is currently being tested on Italian Web pages; however, parts of it should also have wider implications and may be applied to other languages.

The grid includes four layers of description directly related to the writing process. The first of these layers, “Time allowed for writing”, is constructed upon four main categories, related to specific types of texts:

- fast unedited writing (forum postings and, occasionally, Web sites); includes text written without planning and without a second reading and/or correction.
- fast revised writing (forum postings and,

occasionally, Web sites); includes text written with some degree of planning and/correction.

- conventional revised writing (Web sites and, occasionally, forum postings); includes text written within a process of planning and correction.
- writing designed for other kinds of publishing (Web sites and, occasionally, forum postings); includes text written for other media mechanically copied and published on the Web.

What kind of place can be allocated to blogs in this classification? A preliminary answer to this question will be provided in the Conclusion.

Method of inquiry

The editing tendency of a single blogger is, of course, strictly personal and individual blogs show huge variations in linguistic solutions. We can however try to individuate the slant of the textual genre by looking at a great mass of data. The search engine querying of entire blog sites, such as Blog.excite.it, is the simplest way of doing this.

Unedited writing should become apparent at many different levels: lexicon, semantics, syntax and so forth. It is in the orthography, however, that bad editing should be most evident. The frequency of writing errors such as misspellings should be a good index of unedited writing.

We will therefore try to determine the correctness of blogs by measuring the proportion of correct and wrong forms of challenging Italian words. This ratio will be measured in three different situations: the whole Web, some blog sites and some newspaper sites. We must also remember that a text published on the Web may have been written on a word processor with orthographic correction. This kind of tool would easily reduce the number of writing errors. This should happen often in newspaper texts but only occasionally in blogs, where text is probably often entered through web interfaces for direct composition lacking of orthographic correction tools.

As for the errors themselves, the most widely used single-volume Italian dictionary, the *Zingarelli*, lists in a specific entry 103 “frequent” errors in written and spoken Italian (Zingarelli 2005, *ad v.* “errore”; the list is slightly revised and updated through the yearly reprints of the dictionary).

Most of the entries of the list involve errors of pronunciation and are irrelevant to a discussion of written language. Other entries involve errors in graphical accent which in a blog could be difficult to avoid due to writing tool problems. As a matter of fact, the final aspect of an Italian web text may be influenced by the use of particular tools, as discussed in the second layer of the above-mentioned empirical grid:

- keyboard or interface not suitable for the task
- standard keyboard or interface, suitable for the task
- professional tools

Italian Web publishing, especially if connected to Web forms, often encounters problems not just with text formatting (italics, bold etc.), but also with the correct representation of accented letters. Unwanted substitutions of characters are frequent: a writer may type an orthographically correct text only to discover that the publishing system being used cannot handle accented letters or text formatting. The Italian forum *La meglio gioventù* published in 2004 on the Web site of the newspaper *La repubblica* provides many examples of this. Many Italians living abroad have taken part in the forum; many orthographic errors can therefore be explained by the use of non-Italian keyboards (e.g. keyboards without direct access to accented characters) and not by the writers’ lack of orthographic competence. The following is a quotation from a posting from England, where accented letters are replaced by the sequence letter + apex:

Non ho potuto vedere il film perche' [= perché] non ho accesso ai canali RAI in questi giorni e la cosa mi rattrista molto. Penso che la mia meglio gioventu' [= gioventù] sia legata al momento in cui ho cominciato a decidere da sola.

For unequivocal use on the Web, excluding words using accents and apostrophes, the *Zingarelli* list should therefore be reduced to 21 simple errors. We have to reduce it further because the huge numbers involved in search engine queries make it impractical to distinguish homographs. We will therefore avoid discussing the errors *avvallo* for *avallo* and *Macchiavelli* for *Machiavelli*, since the regular word *avvallo* (from the verb *avvallare*) and a common surname *Macchiavelli* do exist in Italian. Results of the selection are shown in Table 1.

We should also note that some of the “wrong” forms in this list are included in standard Italian by other lexicographers. The Italian dictionary of Tullio De Mauro (2000) acknowledges as regular forms *efficenza* (“bureaucratic use”) and *peronospora* (“variant”).

Moreover, this short list of errors is aimed at an high-level writing. Some of the misspellings included in it are therefore likely to be made only by experienced writers: rare words such as *collutazione* or *collutorio* may only be used, even if misspelled, by people with a good command of the Italian language. Errors in commonly used words are instead typical of less experienced writers: in our list, *eccezzionale* or *scenza* seem the most conspicuous examples of this. The list does not include more basic writing errors in commonplace words, such as *propio* instead of *proprio* (see below). It includes instead many words of a technical or literary nature. Anyway, as for the forms selected in Table 1, De Mauro (2000) marks only *collutorio* and *peronospora* as words used only in technical or scientific texts; the remaining words are considered of widespread use or knowledge.

Searching for errors using search engines

Commercial search engines lack many features typical of corpus querying systems. But even their unsophisticated linguistic functions can still be exploited in a significant way (Kilgariff and Grefenstette 2003:342; Calishain and Dornfest 2003; Maxwell 2004; Davis 2005). The trickiest problem, in this field, is probably the fact that the most efficient search engine, Google, provides often only the approximate number of pages where a given token occurs, instead of the exact number of occurrences of the token. Due to this engine behavior, values and figures can only be compared with other search engine data and cannot provide reliable absolute values.

The searches presented below were done using Google in October 2005 – January 2006. All of the searches were restricted to the pages in Italian (using the option of language restriction offered by the engine); the figures provided refer to the “number of pages” or to the approximate number of occurrences found by the search engine. Moreover, we will not take into account conjugated forms or differences between singular and plural, or masculine and feminine.

From the point of view of corpus extension we should note that Google itself enables to search the “whole web” but does not allow restricted searches like “all forums”, “all newspapers” and so on (Google Blogsearch being still in beta testing). Searches were then restricted to single sites through the “site” function of the search engine. We should also note that not all blogs or newspapers, nor the “deep web”, can be accessed using search engines. For example, the popular Italian blogging site Digilander is not indexed by Google and the contents of its blogs cannot then be retrieved in this way. The selection of blog and newspaper sites was then created by trial on the most popular Italian sites of their kind.

The blog sites selected were: Blog.excite.it;

Clarence.com; Splinder.it (contents are partially repeated in the following); Splinder.com.

The newspaper sites selected were: Corriere.it (*Corriere della sera*); Ilmattino.caltanet.it (*Il mattino*); Repubblica.it (*La repubblica*); Unita.it (*L'Unità*).

Search Results

The search results are summarized in Table 1. Figures display the percentage of wrong forms against correct ones (total of correct forms = 100).

Table 1

Wrong form	Correct form	Whole Web	Blog sites	Newspaper sites
accelerare	accelerare	5.83	18.26	4.06
anedottico	aneddotico	205.26	10.61	10.00
appropriato	appropriato	0.06	0.71	0.59
areoport	aeroporto	10.26	1.57	0.37
Caltanissetta	Caltanissetta	18.42	10.57	4.99
collutazione	colluttazione	3.86	29.90	9.50
collutorio	collutorio	31.90	188.74	61.54
conoscenza	conoscenza	0.61	0.22	0.29
coscenza	coscienza	2.15	0.51	0.96
eccezzionale	eccezionale	4.33	0.97	2.13
efficenza	efficienza	3.82	2.52	5.20
essicare	essicare	2.78	10.53	8.06
esterefatto	esterrefatto	25.26	44.60	22.22
ingegnere	ingegnere	0.03	0.30	1.49
Missisipi	Mississippi	0.20	7.60	3.75
metereologia	meteorologia	14.97	56.16	5.26
peronospera	peronospora	2.86	130.77	85.71
rindondante	ridondante	0.32	5.37	1.96
scenza	scienza	1.38	0.11	0.13
Totals		4.28	0.74	0.68

Looking at the results for the whole Web we can see that:

- The balance of wrong / correct forms is highly differentiated. In the case of *anedottico* / *aneddotico* the wrong form is twice as common as the correct one. In the case of *ingegnere* / *ingegnere* the wrong form accounts for only 0.03% of the use of the word.
- The frequency of words is also highly differentiated: as for correct forms, it ranges from 6,480,000 occurrences of *scienza* to 114 of *aneddotico*; as for wrong forms, it ranges from 488,000 of *Caltanissetta* to 234 of *aneddotico*.

Restricting the search to different kind of sites we can then see that blogs show a much lower percentage of errors than the whole Web. In fact, the percentages are often similar to those of professionally edited texts such as Web

newspapers. Both kinds of sites are also consistently more accurate in their spelling than the average of the Italian Web pages. Queries executed through the AltaVista engine do not confirm this result, but blog indexing by this engine seems particularly poor (yielding 1/30 of the Google occurrences, while the ratio for newspapers is nearly 1:1).

In general, we should also note that blogs have a significantly higher percentage of errors than newspapers as for rare or technical words (e.g. *collutorio*). On the other hand they have consistently lower percentages of errors in more common words like *eccezionale* or *coscienza*. However, a 19-word vocabulary is too small to draw conclusions from this observation. Other common errors are in fact more frequent in blogs than in newspapers. Searching for the misspelled word *propio* for *proprio* on the same blog sites of the queries discussed above, we obtain an error percentage of 0.59 (with a maximum of 1.83 and a minimum of 0.31). Testing it with the pre-examined newspaper sites yields an error percentage of only 0.1. The percentage of errors on the entire Web is on the contrary considerably higher than both figures (4%).

Conclusion

The quantity of data and the quality of the sample seem sufficient to draw a preliminary conclusion: on the orthographic level, Italian blogs are edited with the same care as Italian online newspapers. There are of course differences but on average a rough index such as the total percentages of errors taken from a limited list gives blogs and newspapers equal rank while distancing them from the less edited mass of Web texts. Pending further analysis, we can then tentatively include the average of those blogs in the category of “conventional revised writing”.

References

- Calishain, Tara, and Dornfest, Rael. 2003. *Google Hacks*. Beijing, etc.: O'Reilly.
- Crystal, D. 2001. *Language and the Internet*. Cambridge: Cambridge University Press.
- Davis, H. 2005. *Building Research Tools With Google for Dummies*. Hoboken: Wiley Publishing.
- De Mauro, T. 2000. *Il dizionario della lingua italiana*. Torino: Paravia.
- Kilgariff, A. and Grefenstette, G. 2003. Introduction to the Special Issue on the Web as Corpus. *Computational Linguistics* 29(3):333-347.
- Maxwell, M. 2004. Resource Discovery for Low Density Languages: Internet Internet Search, abstract in the abstract book of ACH/ALLC 2004 - Goteborg University, Goteborg: 88-89.
- Tavosanis, M. 2005. *Linguistic Variability of Web Italian: a Working Empirical Grid*. Forthcoming (included in the ACH/ALLC 2005 program – Victoria University).
- Zingarelli 2006. 2005. Bologna: Zanichelli.