

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
PIHHI (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

SOI: [1.1/TAS](#) DOI: [10.15863/TAS](#)

International Scientific Journal Theoretical & Applied Science

p-ISSN: 2308-4944 (print) e-ISSN: 2409-0085 (online)

Year: 2019 Issue: 10 Volume: 78

Published: 14.10.2019 <http://T-Science.org>

QR – Issue



QR – Article



Nozimjon Bobojon o'g'li Ataboev

Uzbek State World Languages University

Ph.D. student,

Tashkent, Uzbekistan

+998919254881

anb929292@gmail.com

PROBLEMATIC ISSUES OF CORPUS ANALYSIS AND ITS SHORTCOMINGS

Abstract: The article deals with the problems of corpus linguistics and corpus analyses. The use of corpus results can be problematic as well as the shortcomings might emerge while applying corpus analyses in the different linguistic fields. The practical analyses to prove the dysfunctions of the corpus application in terms of searches as well as quantitative results have been undertaken. The theoretically valuable data included in article were based on the works by Vsevolodova M., Makarov M., Magomedova A., Kopotev M. and Fillmore Ch. The corpus-based quantitative results derived from the modern corpus 'COCA' can increase the practical value of the work. The final conclusions have been made after undertaking several searches on the concordance.

Key words: Corpus linguistics, dysfunctions of corpora, corpus, corpus analyses, concordance, collocations.

Language: English

Citation: Ataboev, N. B. (2019). Problematic issues of corpus analysis and its shortcomings. *ISJ Theoretical & Applied Science*, 10 (78), 170-173.

Soi: <http://s-o-i.org/1.1/TAS-10-78-30> **Doi:**  <https://dx.doi.org/10.15863/TAS.2019.10.78.30>

Scopus ASCC: 1203.

Introduction

New information technologies lead to the opportunity to learn the language not only from the traditional "storages" of language data, such as dictionaries, artworks, classics, written texts as representative samples from the general population, but also to enter into the computer and process large collected texts – CORPUS. In this regard, we have a qualitative leap, which is especially often noted in lexicology and lexicography: for example, if you used to compile dictionaries, which is an extremely important theoretical and practical work that determines, in essence, the composition and structure of all linguistic research and the application of linguistics to practical problems, so important, for example, as teaching language and speech, now this work is much easier. The possibilities for observing and studying speech, both oral and written, have expanded dramatically. The huge language material that modern computers are able to process makes it possible to test the proposed theoretical models of linguistic phenomena and develop new ones [7, p. 2].

Many traditional problems of linguistics are being solved in a new way and their solution to traditional problems is achieved much easier and more conclusively.

The field of corpus linguistics (CL), as well as projects of electronic corpus of texts are actively developing and occupying the leading positions in the methodology of teaching foreign languages, having significant applied potential. In the process of teaching foreign languages, one of the main problems is the lack of appropriate pedagogical textual materials and relevant grammatical models. The linguistic corpus is considered as one of the modern information resources, on the basis of which it is possible to form the lexical and grammatical speech skills of students.

Discussion

The use of the corpus facilitates the search for material, reduces the complexity, reduces the time spent, provides high accuracy of the selection and the reliability of the research results. The corpus, in fact, is a continuation of the file cabinets that linguists have always worked with, however, according to scholars,

Impact Factor:

ISRA (India)	= 4.971	SIS (USA)	= 0.912	ICV (Poland)	= 6.630
ISI (Dubai, UAE)	= 0.829	PIHHI (Russia)	= 0.126	PIF (India)	= 1.940
GIF (Australia)	= 0.564	ESJI (KZ)	= 8.716	IBI (India)	= 4.260
JIF	= 1.500	SJIF (Morocco)	= 5.667	OAJI (USA)	= 0.350

it allows one to "get out of the analysis of individual" correct "sentences. in a sea of real natural texts in a machine-readable format ". Many authors, however, come to the conclusion that in order to obtain reliable results, the complementary use of both corpus data and traditional sources of material is necessary [2].

The corpus data, in contrast to the data of dictionaries and grammar guides subjected to filtering by people with a deep knowledge of the language and linguistic instinct, give an idea not about the standard of the language, but about the whole diversity of its individual, dialect and genre variation. Due to this, the study of corpus data, writes G. Lorenz, allows a less rigid approach to what is considered to be a violation of the rules of the language norm [Lorenz 2001]. In general, the linguistic corpus "crystallizes problem areas" in the description of the language, being a source for improving theoretical approaches to it [6, p. 151].

The fact that the language corpus contains texts of different genres is also its great advantage as a source of material over the combination of texts of fiction traditionally used by researchers. The latter, according to M. V. Vsevolodova, is "even the most perfect - but only one of the functional styles" [12, p. 131]. The national corpus includes samples of both fiction and texts of other functional styles, which makes it possible to get an idea of modern speech in all its diversity. Of particular note is the importance of the presence of media texts in the corpus material, which in our time, according to the correct remark of V. I. Makarov, perform a normalizing function [8].

Results

The term "corpus" is often used and misinterpreted, coexisting and often merging with such concepts as "collection of texts", "full-text database", "electronic archive", "electronic library" [11][13]. In a broad sense, the corpus refers to any combination of texts connected by one characteristic - authorship, genre, etc. [5]. Corpus also includes datasets, which are essentially thesauruses, or meta-corpus. Recently, there has been a tendency to consider as the body of texts and the entire information space of the Internet [13] [4] [3] [10].

In the narrow sense, the corpus of texts (CT) is understood as a unified, structured and labeled collection of language (speech) data in electronic form [13]. The definition of the corpus adopted in corpus linguistics is based on four main features: 1) the location of the corpus on a machine medium, 2) a standardized representation of the verbal material on this medium, 3) the final size, 4) representativeness as a result of a special selection procedure [11] [9]. The most significant feature is the representativeness, which, in essence, distinguishes CT in a narrow sense.

As it is clear, corpus linguistics is a new field of applied linguistics and being a modern as well as currently developing sphere, it has some disadvantages and several theoretical problems.

To inform about some problematic features of CL, its *disadvantageous* sides have been analyzed. It is evidence that *Corpus linguistics is not able to:*

✓ *provide negative evidence:* this means that a corpus cannot tell us what is possible or correct or not possible or incorrect in language; it can only inform us what is and is not present in the corpus.

✓ *explain why:* CL cannot explain why something is the way it is, only tell us what it is. To find out why, we, as users of language, use our intuition.

✓ *provide all the possible language at the one time:* By the definition, a corpus should be principled: "a large, *principled* collection of naturally occurring texts..." meaning that the language that goes into a corpus is not random, but planned. However, no matter how planned, principled, or large a corpus is, it cannot be a representative of a language.

It is time to discuss the next problem for CL. This is the problem of *authenticity* in the language data supplied by corpora. It is often argued that corpora provide learners with 'authentic' or 'real' language, and since these words echo the key features of Communicative Language Teaching (CLT, hereafter) method that favors the use of authentic and real language over concocted ones, it is often assumed that corpus-based language materials are well-suited for CLT. However, some of the researchers have cast doubt on whether language data in corpora are truly authentic. Widdowson contrasted the concept of 'genuineness' and 'authenticity' and argued that 'genuineness' is the property of texts and is an absolute quality, while 'authenticity' is the characteristic of discourse interpretation. He claimed that language in corpora can be genuine, but it is not authentic because it is isolated from discursive and communicative nature of language.

Moreover, it would be relevant to give some information about the following challenging issue. That is how *to measure the proportion* that dialogs make up of the speech of one particular group, for example, adolescents. Corpus compilers can only record a tiny sample of all adolescents, and how would they measure the proportion of dialogs – in terms of time? in terms of sentences? in terms of words? And if they tried to compile a corpus representative of a language as a whole, then how would they measure the importance of a particular linguistic variety? As we can see that a corpus is not always a reliable database of a language or a sublanguage in terms of the mentioned problematic items.

In addition to that, one of the biggest dysfunctions of corpora can be seen from the following quote: "I don't think there can be any

Impact Factor:

ISRA (India) = 4.971	SIS (USA) = 0.912	ICV (Poland) = 6.630
ISI (Dubai, UAE) = 0.829	PIHHI (Russia) = 0.126	PIF (India) = 1.940
GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

corpora, however large, that contain information about all of the areas of English lexicon and grammar that I want to explore”. [2, p. 35] It is of course true that the sheer *volume of natural language* will never be able to be *captured* inside a database because it is truly mathematically infinite.

Now some examples can be indicated about the use of corpus analyses and their shortcomings. For

example, if the user wants the corpus data regarding use of a word or set of words, he/she needs to know what exactly to be searched. That means, a corpus is an electronic data which is, in its terms, not equipped with a virtual mind to select the words for them.

The searches have been done in the corpus of COCA (Corpus of Contemporary American English). The following data has been taken in the table 1.

Table 1. Collocational unit(CU)s preceding ‘Mother’

CUs	Freq. in CUs	Total frequency
Single	1449	89252 – 1.62 %
Birth	690	32265 – 2.14 %
Teresa	686	4259 – 16.19 %
Biological	324	17799 – 1.82 %
Holy	162	18697 – 0.87 %
Adoptive	205	1894 – 10.82 %
Unmarried	32	2206 – 1.45 %
Surrogate	171	2118 – 8.07 %

(The table was generated from the data derived from COCA [1])

Analysing the corpus-based results, it has been a clear fact that a language corpus is not always something that can provide with the valuable and applicable information. As it has been mentioned above, the given examples are not defined or explained. The searches on concordance have been undertaken for the collocational units preceding the word ‘mother’ and quantitative results have been taken, among which some of them might seem unknown. In order to clarify the meaning of the given CUs, one needs to take further steps because a corpus does not provide any evidence regarding the meaning or the usage of the words.

For example, the CU of Teresa Mother was used many times, i.e. 686 times that is the 16.19 % - in almost a fifth of the total frequency (4259 times) of the ‘Teresa’ was collocated with the token ‘mother’. Of course, it has drawn our attention and we wanted to analyze the meaning of the CU phrase. When the contextual and wider contextual format of the use of the CU has shown the following pieces of the text:

1. “I already respected you immensely, but even more so. You're like **Mother Teresa** to me. We'll be right back. Thank you. Welcome back...” [1]

2. “**Mother Teresa's** sainthood gains broader meaning in view of her' dark night of the...”[1]

3. “Canonization recognizes holiness, not perfection. In elevating **Mother Teresa** to the honors of the altar, the Roman Catholic Church is not overlooking...” [1]

After seeing the above-mentioned data analyses, one can make a conclusion that ‘Mother Teresa’ is a symbol of holiness and is one of the religious notions. As can be seen, the data given in the corpus can be enough to make an assumption but not an exact interpretation or definition for the searched token.

Conclusion

To sum up, it would be essential to note that every science in its emergence experiences some problems. In fact, Corpus linguistics also has several challenges as mentioned above. However, those have had no proper solutions yet. As new investigators in CL, we believe that there will be undertaken enough researches in order to sort the problems out.

Impact Factor:	ISRA (India) = 4.971	SIS (USA) = 0.912	ICV (Poland) = 6.630
	ISI (Dubai, UAE) = 0.829	PIHHI (Russia) = 0.126	PIF (India) = 1.940
	GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
	JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

References:

- (n.d.). *Corpus of Contemporary American English*. Retrieved September 25, 2019, from <https://www.english-corpora.org/coca/>
- Fillmore, C. (1992). "Corpus linguistics or Computer-aided armchair linguistics." J. Svartvik (ud.) *Directions in Corpus Linguistics*. Proceedings of Nobel Symposium 82, pp. 35–60.
- Hemming, C., & Lassi, M. (n.d.). *Copyright and the Web as Corpus Text*. Retrieved 2019, from www.hemming.se/gslt/copyrightHemmingLassi.pdf
- Kilgariff, A. (2001). *Web as corpus Text*. A. Kilgariff // Proc. of Corpus Linguistics 2001 conference (Lancaster University). (pp.342-344). Lancaster.
- Klimenko, S. V., & Rykov, V.V. (2000). *Korpus tekstov kak princip samoorganizacii predmetnoj oblasti Tekst. (Korpus tekstov kak princip samoorganizacii predmetnoj oblasti Tekst.)* Trudy konferencii Dialog. Jelektron, izd. - Retrieved 2019, from <http://rykov-cl.narod.ru/t.html>.
- Kopotev, M.V., & Janda, L. (2006). Nacional'nyj korpus ruskogo jazyka.(Recenzija) Tekst. [National Corpus of the Russian Language. (Review) Text.] *Voprosy jazykoznanija*, №5, pp. 149-155.
- Magomedova, A. N. (2004). *Korpusnaja lingvistika i kontekstnoe razreshenie leksicheskoj mnogoznachnosti [Corpus linguistics and contextual resolution of lexical ambiguity]* slov.diss.kandidat filologicheskikh nauk. (p.203). Mahachkala.
- Makarov, M. L. (2003). *Osnovy teorii diskursa [Fundamentals of the theory of discourse.]* (p.280). Moscow: ITDGK «Gnozis».
- McEnery, T., & Wilson, A. (1999). *Corpus Linguistics Text*. Edinburgh: Edinburgh University Press. Retrieved 2019, from http://bowland-files.lancs.ac.uk/monkey/ihe/linguistics/content_s.htm
- Powell, J. (n.d.). *The Web as a Corpus: A Survey Text*. Retrieved 2019, from www.personal.umich.edu/~jcpowell/WebAsCorpus.rtf
- Rykov, V.V. (n.d.). *Korpus tekstov novyj tip slovesnogo edinstva Tekst. (Text Corpus is a new type of verbal unity Text.)* Jelektron, izd. - Rezhim dostupa k izd. Retrieved 2019, from <http://rykov-cl.narod.ru/t31.html>
- Vsevolodova, M.V., & Seliverstova, O.N. (2006). Trudy po semantike. (Recenzija) Tekst. [Works on semantics. (Review) Text.] *Voprosy jazykoznanija*, № 3, pp. 130-135.
- Zaharov, V.P. (2005). *Korpusnaja lingvistika. (Corpus linguistics)* Uchebno-metodicheskoe posobie. SPb.