

In our opinion, sociology can lose essential meanings in the research process because of lack of in-deep immersion in the daily life and speech of the people. Context can be understood as institutional frame, which defines the status and role of the concept in social structure. We chose the concepts of “altruism” and “mercy” as examples to demonstrate the corpus-based conceptualisation and its place in sociological research methodology.

**Maria Rubtsova
Elena Vasilieva
Oleg Pavenkov
Vladimir Pavenkov**

Corpus-based conceptualization in sociology: possibilities and limits

Conceituação corpus-based em sociologia: possibilidades e limites

MARIA RUBTSOVA*

ELENA VASILIEVA**

OLEG PAVENKOV***

VLADIMIR PAVENKOV****

Abstract

The problem of the quantitative interpretation of qualitative data is one of the most important in sociological research. Textual analysis has placed emphasis on deep and careful study of texts how personal strategies embodied in the concepts. However, quantitative interpretation has always been problematic. Our paper deals with the corpus-based conceptualization method, which can be considered as a method of collecting and organizing data material from linguistic corpora. The corpus-based conceptualization allows us to establish a closer link with the meaning and identify the whole spectrum of meanings. It shows that some sociologists lose essential meanings in the research process because of lack of in-deep immersion in the daily life and speech of the people. We chose the concepts of “altruism” and “mercy” as examples to demonstrate the corpus-based conceptualization and its place in sociological research methodology. Data comes from the Russian National Corpus. The Russian National Corpus consists of 1802 relevant words, with 775 for altruism and 1047 for mercy. Data processing carried out by SPSS 19.0. As the result, we have discussed what difficulties the researcher can meet using this method and have offered SFL and Role and Reference grammar as a way to accurately determine the context. Our

* Doctoral Degree in Herzen State University; Associate Professor in the Department of Social Management and Planning - Saint-Petersburg State University, Russia; Email: infosoc@bk.ru

** Candidate Degree in St.-Petersburg State University, Russia; Leading Researcher - Academy of Sciences of Republic of Sakha (Yakutia), Russia; Email: dsndfn@yandex.ru

*** Candidate (PhD) Degree in Philosophy in Leningrad State University, Russia; Associated Professor in the Department of Media communication - Saint Petersburg Institute of Film and Television, Russia; Email: pavenkov@yahoo.es

**** Candidate (PhD) Degree in History in Leningrad State University, Russia. Associated Professor - Department of Media communication at the Saint Petersburg Institute of Film and Television, Russia; Email: infosoc1@yandex.ru

suggestions can be used in the preparation of questionnaires, guides, in an analysis of interview transcripts.

Keywords: Corpus-based conceptualization. Sociological operationalization. Corpus linguistics.

Resumo

O problema da interpretação quantitativa de dados qualitativos é um dos mais importantes na pesquisa sociológica. A análise textual colocou ênfase no estudo profundo e cuidadoso de textos como estratégias pessoais incorporadas nos conceitos. No entanto, a interpretação quantitativa sempre foi problemática. Nosso artigo trata do método de conceituação *Corpus-based*, que pode ser considerado como um método de coleta e organização de material de dados a partir de corpora linguística. A conceituação *Corpus-based* nos permite estabelecer uma ligação mais próxima com o significado e identificar todo o espectro de significados. Isso mostra que alguns sociólogos perdem significados essenciais no processo de pesquisa por causa da falta de imersão profunda na vida cotidiana e na fala das pessoas. Escolhemos os conceitos de “altruísmo” e “misericórdia” como exemplos para demonstrar a conceituação *Corpus-based* e seu lugar na metodologia de pesquisa sociológica. Os dados vêm do Corpus Nacional Russo. O Corpus Nacional Russo consiste em 1802 palavras relevantes, com 775 para o altruísmo e 1047 para a misericórdia. Processamento de dados realizado pelo SPSS 19.0. Como resultado, discutimos as dificuldades que o pesquisador pode encontrar usando esse método e oferecemos a gramática *SFL and Role and Reference* como uma maneira de determinar com precisão o contexto. Nossas sugestões podem ser utilizadas na elaboração de questionários, guias, em uma análise de transcrições de entrevistas.

Palavras-chave: conceituação *Corpus-based*. Operacionalização sociológica. Lingüística de Corpus.

Introduction

The problem of the quantitative interpretation of qualitative data is one of the most important in sociological research. Textual analysis has placed emphasis on deep and careful study of texts how personal strategies embodied in the concepts. However, quantitative interpretation has always been problematic.

Quantitative operationalization of concepts can improve the objectivity of the data and exclude some erroneous or obscure use of social categories, as well as determine the particular application in different contexts to identify the relevance of social problems. As stated by P. Lazarsfeld, operationalisation of concepts can help to create a model that characterizes the social problems studied in the course of empirical social research [Lazarsfeld, 1962]. This

model is created on the stage of conceptualization of the concept and is used as a base for research operations.

We propose to use quantitative analysis as the basis for building empirical models, which allows identifying the meanings and trends in use through the use of corpus linguistics. Research and institutional lexical scoping to determine the spread of certain values in the various subjects, which allows on the basis of a formalized conceptualization of the operationalization of concepts within the framework of sociological research and to increase the objectivity of the research conducted.

Our paper deals with the corpus-based conceptualization method, which can consider as a method of organizing data using linguistic corpora. The corpus-based conceptualization allows us to establish a closer link with the meaning and identify the whole spectrum of meanings.

First of all, corpus linguistics studies a language not only as a remarkable social phenomenon (e.g. altruism, power, government). The most significant aim is to explain the language characteristics of a term (concept) (Corpora and Discourse, 2008). Corpus does not initially have such intention; however, it creates for sociology some interesting content that extends standard sociological frames of social categories' knowledge (Abulof, 2015). The attention of the linguistic corpus in natural language gives a path to answers of how people actually discuss social reality (Blei, 2012). If standard sociological content analysis requires conceptual study through the subjective descriptions of a sociologist, the material is already presented in a corpus, and sociologists can only use what is available. This increases the level of objectivity in a researcher's operationalisation.

Second, a corpus has a meaningful size that allows us to make a big data statistical analysis. The most part of the national corpora contain more than 500 million words (Corpora & Discourse, 2008). Materials are described by years and months over a long period; for example National Corpus of the Russian Language includes the development of language from the 18th to the beginning of the 21st century (National Corpus of the Russian Language, 2015). A researcher can build the frequency of the use of categories over various periods of time (diachronic word frequency): e.g., the max. frequency of use of the concept 'mercy' is in years 2002-2004. The aggregate of phrases also supports the in-depth study of the dynamic context of the concept and its main topics (see, e.g.: Liu, 2012; Abulof, 2015).

Also the main part of national corpora have one principle of structuring that is crucial for cross-cultural studies (Corpora & Cross-linguistic Research, 1998). The corpus linguistics gives us the possibility to recognize differences in the understanding of the social terms in different cultures that should be done before comparative sociological studies (Rawoens, 2010).

Finally, the corpus research is not expensive and time-consuming. It can be done in a user-friendly search format for uploading and downloading data that is compatible with Word, Excel, SPSS etc. (McEney & Hardie,

2012). So we can do the interpretation procedure quickly. Researchers are excited by the corpus linguistics possibilities because it allows to analyze a huge amount of text inexpensively (Schonhardt-Bailey, 2005)

In our opinion, sociology can lose essential meanings in the research process because of lack of in-deep immersion in the daily life and speech of the people. Context can be understood as institutional frame, which defines the status and role of the concept in social structure. We chose the concepts of "altruism" and "mercy" as examples to demonstrate the corpus-based conceptualisation and its place in sociological research methodology. An "altruism" is defines as the willingness to act selflessly for the benefit of others, without regard to their interests. A "mercy" – as the willingness to help someone or to forgive someone out of compassion, humanity.

Specifically, the current study attempts to answer the following research questions:

- Are the concepts «altruism» and «mercy» interchangeable as synonyms in sociological research?

Data and methodology

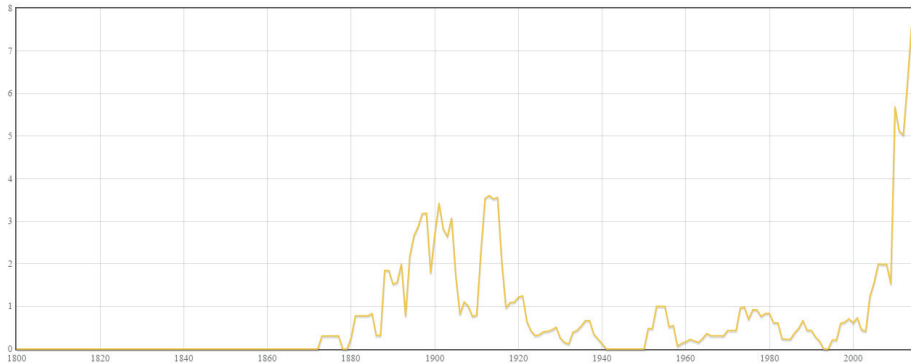
Research design

The proposed methodology of the study of social categories is the combination of quantitative and qualitative analysis of the encoded array of a Corpus. Quantifying the frequencies of word used was based on expert's context encoding. In accordance with the method, three independent experts performed the encoding. All of them were representatives of St. Petersburg universities and were not the article's authors. They offered the following seven contexts: people's actions, humans themselves, quality of relations, state, social institutes, organisations, conception-ideology. Data processing was carried out using SPSS.

Sampling procedures

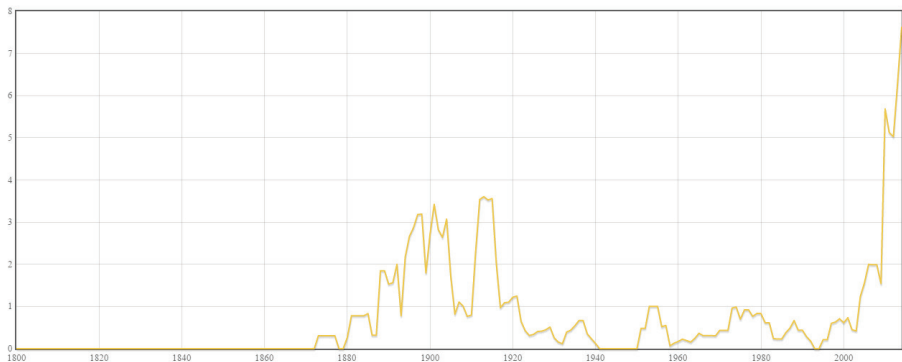
Data comes from the Russian National Corpus. In Fig.1, 2 we can see the frequency of references to the category of "mercy" in the period of 1800-2014 years. The Russian National Corpus consists of 1802 relevant words, with 775 for altruism and 1047 for mercy. Data processing carried out by SPSS 19.0. Because Russian has two concepts denoting the social phenomena - 'altruism' and 'mercy' (mutual help), which are considered as synonymous, we will explain their most frequent use in the main, newspaper and spoken Russian Corpora (Russian National Corpus, 2016).

Fig. 1 - The frequency of references to the category of "mercy" in the period of 1800-2014 years, the number of references per year



Source: Russian National Corpus URL.: <http://www.ruscorpora.ru/>

Fig. 2 - The frequency of references to the category of "altruism" in the period of 1800-2014 years, the number of references per year



Source: Russian National Corpus URL.: <http://www.ruscorpora.ru/>

Results

Both categories are mentioned in the following institutional contexts: "People`s actions", "Human himself", "Relations` quality", "State", "Social institutes", "Organisations", "Conceptions and ideology". The studying of the words "altruism" and "mercy", according to these contexts, led to the results discussed below.

As we can see in fig.1 and fig.2 the concept "altruism" is less popular in Russian language, then the concept "mercy": there is 1047 mentions of the mercy and only 309 mentions of altruism in corpora. The statistical analysis showed that the differences in use of these concepts have statistical significance: Pearson Chi-Square $p < 0,000$ (see Table 1).

Table 1 - Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	424,507a	6	,000
Likelihood Ratio	409,543	6	,000
Linear-by-Linear Association	252,619	1	,000
N of Valid Cases	1356		

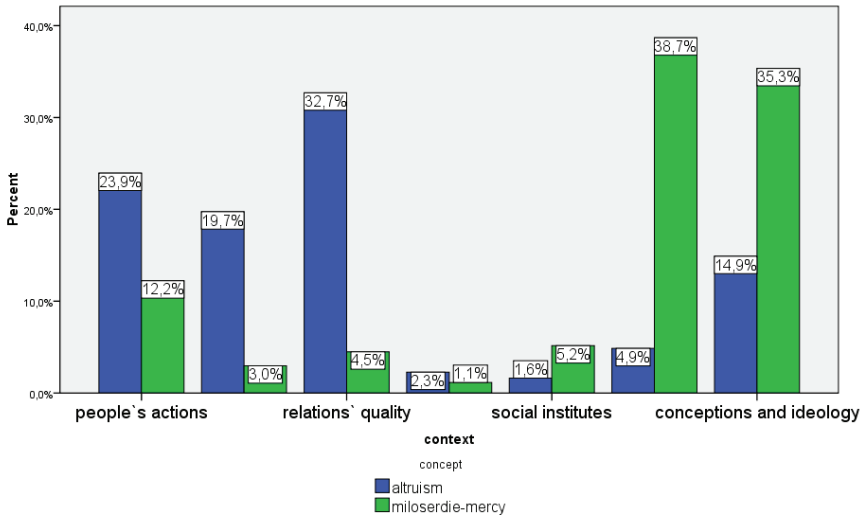
a) 1 cells (7,1%) have expected count less than 5. The minimum expected count is 4,33.

The analysis of mentioning of the words “altruism” and “mercy” in distinct contexts shows, that these concepts has different semantic load. The altruism more often use in context “relations` quality” and “people`s actions”, i.e. as personality characteristic; in opposition, the mercy is used in context “organisations” and “conceptions and ideology”, i.e. as social feature (see Table 2 and Figure 3).

Table 2 - The frequency of references to the category of “altruism” and “mercy” in contexts

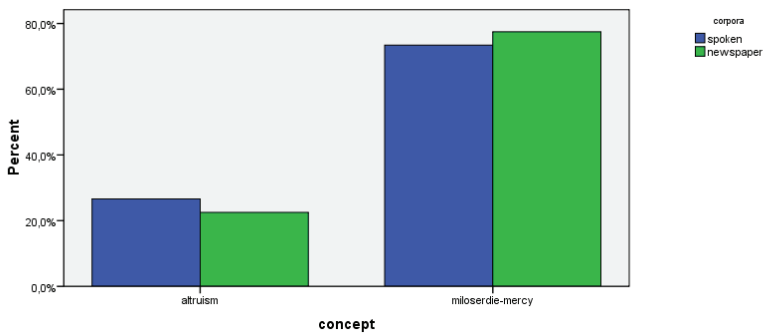
			concept		Total
			altruism	mercy	
context	people`s actions	Count	74	128	202
		% within concept	23,9%	12,2%	14,9%
	human himself	Count	61	31	92
		% within concept	19,7%	3,0%	6,8%
	relations` quality	Count	101	47	148
		% within concept	32,7%	4,5%	10,9%
	state	Count	7	12	19
		% within concept	2,3%	1,1%	1,4%
	social institutes	Count	5	54	59
		% within concept	1,6%	5,2%	4,4%
	organisations	Count	15	405	420
		% within concept	4,9%	38,7%	31,0%
	conceptions and ideology	Count	46	370	416
		% within concept	14,9%	35,3%	30,7%
Total		Count	309	1047	1356
		% within concept	100,0%	100,0%	100,0%

Fig. 3 - The frequency of references to the category of "altruism" and "mercy" in different contexts (%)



Also there is different of using of analyzed concepts in spoken corpora and in newspaper: the altruism more typical for speeches, while mercy more frequently used in written sources (see Figure 4).

Fig. 4 - The frequency of references to the category of "altruism" and "mercy" in spoken corpora and in newspaper (%)



The altruism in spoken corpora is used in context "conceptions and ideology" (56% of all mentions in spoken corpora) and "human himself" (44%), while in newspapers it increasingly used in context "relations' quality" (35.6%) and "people's actions" (26.1%). This could means that people treat this concept as a part of ideals, beliefs and principles of social life, but in official sources it is handled as an impulsive cause.

The mercy in spoken corpora is also used in context "conceptions and

ideology" (58%), but in newspapers it is frequently mentioned in context "organisations" (40%). It possibly means, that this word in written sources is used as reflection of action of charities. In the same time, a significant level of mentions of mercy in newspapers in context "conceptions and ideology" (40%), shows that this concept has big value for morality and ethical basis of society (see Table 3).

Table 3 - The frequency of references to the category of "altruism" and "mercy" in different corpora

concept		corpora					
		spoken	newspaper	Total			
altruism	context	people`s actions	Frequency % in corpora	0 ,0%	74 26,1%	74 23,9%	
		human himself	Frequency % in corpora	11 44,0%	50 17,6%	61 19,7%	
	relations` quality	Frequency % in corpora	0 ,0%	101 35,6%	101 32,7%		
	state	Frequency % in corpora	0 ,0%	7 2,5%	7 2,3%		
	social institutes	Frequency % in corpora	0 ,0%	5 1,8%	5 1,6%		
	organisations	Frequency % in corpora	0 ,0%	15 5,3%	15 4,9%		
	conceptions and ideology	Frequency % in corpora	14 56,0%	32 11,3%	46 14,9%		
	Total	Frequency % in corpora	25 100,0%	284 100,0%	309 100,0%		
	mercy	context	people`s actions	Frequency % in corpora	2 2,9%	126 12,9%	128 12,2%
			human himself	Frequency % in corpora	5 7,2%	26 2,7%	31 3,0%
relations` quality		Frequency % in corpora	4 5,8%	43 4,4%	47 4,5%		
state		Frequency % in corpora	1 1,4%	11 1,1%	12 1,1%		
social institutes		Frequency % in corpora	3 4,3%	51 5,2%	54 5,2%		
organisations		Frequency % in corpora	14 20,3%	391 40,0%	405 38,7%		
conceptions and ideology		Frequency % in corpora	40 58,0%	330 33,7%	370 35,3%		
Total		Frequency % in corpora	69 100,0%	978 100,0%	1047 100,0%		

Discussion and conclusion

We explored the everyday using of the two Russian words, 'altruism'

and 'mercy', which originally had the same meanings. The concepts can be considered as synonymous in questionnaires and guides of interview. Based on an analysis by the Russian National Corpus, we have described seven contexts of word use for 'altruism' and 'mercy': the following institutional contexts: "People`s actions", "Human himself", "Relations` quality", "State", "Social institutes", "Organisations", "Conceptions and ideology". The quantitative analysis shows differences in the use of these concepts. The differences between 'mercy' and 'altruism' are statistically significant for all three Russian corpora. Thus, their use as a synonym in sociological research is wrong.

The analysis shows both advantages and disadvantages of the proposed method of corpus linguistics for the operationalization of social categories. The positive side is our possibility to give a description of the context of its usage. We put emphasis on the noun (concept) that avoids the loss of meaning, which the respondents may use in answering the survey questions, and especially in the case of in-depth semi-structured interviews.

However, the proposed method has several disadvantages. Despite the fact that it is aimed at avoiding subjectivity and it is trying to raise objectivity in the operationalization of social categories, the encoding process itself continues to depend on the will of researchers. Existing techniques to reduce subjectivity, such as the coding of several (usually three) independent research has the elements of subjectivity because they are usually the representatives of one culture and affiliated (related) persons. One attempt to avoid subjectivism is an appeal to the systemic functional linguistics (SFL) (M. Halliday) (see Table 4.) and Role and Reference grammar (RRG) (Van Valin)

Table 4 – Comparison: Traditional v. Functional Grammar

Participant (nominal group)	Process (verbal group)		Participant (nominal group)			Circumstance (adjectival group)		
	Noun	Verb (simple past)	Verb (infinite)	Possessive determiner	Adjective	Noun	Preposition	Determiner
Annie	wanted	to play	her	new	clarinet	in	the	band.

Source: Functional Grammar for Teachers

URL:<http://manxman.ch/moodle2/mod/resource/view.php?id=132>

In our example with «altruism» and «mercy», we used conceptualization of a noun, whereas SFL and especially RRG offer us to change the focus on the process (verbal group). Verbal group exists regardless and independent of the investigator and can show the context of the roles of participants. In this case, we can avoid subjectivity in coding and may analyse grammatical structure of sentence. Our suggestions can be used in the preparation of questionnaires, guides, in an analysis of interview transcripts.

References

Abulof, U. (2015). Normative concepts analysis: Unpacking the language of legitimation. *International Journal of Social Research Methodology*, 18(1), 73-89. <http://dx.doi.org/10.1080/13645579.2013.861656>

Aijmer, Karin & Bengt Altenberg (eds.) (1991). *English Corpus Linguistics: Studies in Honour of Jan Svartvik*. London: Longman

Corpora and Discourse: The Challenges of Different Settings (2008). T. 31. P. 1–297. <http://www.corpus-linguistics.com>

Functional Grammar for Teachers (2016).

URL:<http://manxman.ch/moodle2/mod/resource/view.php?id=132> (date treatment: 02/03/2016)

Gladkova, A., & Romero-Trillo, J. (2014). Ain't it beautiful? The conceptualization of beauty from an ethnopragmatic perspective. *Journal of Pragmatics*, 60, 140-159. <http://dx.doi.org/10.1016/j.pragma.2013.11.005>

Halliday, M. A. K. (1975). Sociological aspects of semantic change. In Luigi Heilman (ed), *Proceedings of the Eleventh International Congress of Linguists*. Bologna: Mulino.

Halliday, M. A. K. (1985). *Introduction to Functional Grammar*, London: Edward Arnold.

Halliday, M. A. K. (1999). "Corpus studies and probabilistic grammar". In *English Corpus Linguistics* ed by K. Aijmer & B. Altenberg, 30-43. London: Longman.

Halliday, M. A. K. (1992). "Language as system and language as instance: the corpus as a theoretical construct". Mouton de Gruyter, *Directions in Corpus Linguistics: Proceedings of Nobel Symposium 82, stockholm, 4-8 August 1991*, ed. Jan Svartvik (*Trends in Linguistics Studies and Monographs* 65).

Lazarsfeld, P. F. (1962). *American Sociological Review*, 27(6), 757-767. <http://dx.doi.org/10.2307/2090403>

McEnery T., Hardie A. (2012). *Corpus linguistics: Method, theory and practice*. Cambridge: Cambridge University Press.

McEnery T., Wilson A. (2001). *Corpus Linguistics*. Edinburgh: Edinburgh University Press.

O'Donnell, Michael (1995). "From Corpus to Codings: Semi-Automating the Acquisition of Linguistic Features", in *Proceedings of the AAAI Spring Symposium on Empirical Methods in Discourse Interpretation and Generation*, Stanford University, California, March 27 - 29.

Pavenkov, O.; Pavenkova, M. (2016). Discourse analysis based on Martin and Rose's taxonomy: a case of promoting student discourse on the CLIL PhD programme in religion philosophy. *Reveleto-Revista Electronica Espacio Teologico* 10(17): 129-139.

Pavenkov O.V., Rubtcova M. (2016). Interacción de los géneros académicos Ruso e Inglés en el Aprendizaje Integrado Contenido-Lengua de los programas de doctorado de Sociología de la Administración: una brecha en el proceso de implementación. *Revista Dilemas Contemporáneos: Educación, Política y Valores*. 2016, №1, Enero.

Pavenkov O. (2014). Contemporary linguistic analysis of the concept "love". *Studia Humanitatis*. № 4. – URL: <http://st-hum.ru/content/pavenkov-os-contemporary-linguistic-analysis-concept-love> (Accessed on: 23.07.2016)

Rubtcova, M.; Pavenkov, O.; Khmyrova-Pruel, I.; et al (2016). Systemic functional linguistics (SFL) as sociolinguistic and sociological conception: Possibilities and limits of theoretical framework. *International Journal of Applied Linguistics and English Literature* 5(3): 272-281 DOI: 10.7575/aiac.ijalel, v. 5, n. 3, p. 272.

Rubtcova, M., Pavenkov, O., Pavenkov, V., & Vasilieva, E. (2015). The language of

- altruism: Corpus-based conceptualization of social category for management sociology. *Asian Social Science*, 11(13), 289-297. doi:10.5539/ass.v11n13p289
- Rubtcova, M., Pavenkov, O., Pavenkov, V., & Vasilieva, E. (2015). Representations of trust to public service in Russian newspapers during election time: Corpus-based content analysis in Public administration sociology. *Mediterranean Journal of Social Sciences*, 6(4S1), 436-444. doi:10.5901/mjss.2015.v6n4s1p436
- Rubtcova, M., Pavenkov, O., Varlamova, J., Kaisarova, V., Volchkova, L., Menshikova, G., & Denisova, J. (2016). How to identify negative attitudes towards inclusive education: Critical discourse analysis of Russian transcripts using Role and Reference grammar. *International Journal of Applied Linguistics and English Literature*, 5(5), 183-196. doi:10.7575/aiac.ijalel.v.5n.5p.183
- Rubtsova, M. V., & Vasilieva, E. A. (2016). Conceptualization and operationalization of notion "trust" for applied sociological research. *Sotsiologicheskie Issledovaniya*, 2016(1), 58-65.
- Rubtcova, M. V. (2015). Corpus Linguistics in Sociological Research. *Philosophy of Science eJournal*, 8(11), March 26. <http://sociology.management/2015/03/12/corpus-linguistics-in-sociological-research/>
- Rubtcova M., Pavenkov O. Systemic Functional Linguistics as a Macro-sociolinguistics Framework: The Stages of Development. // The Joongwon Linguistic Society of Korea. *Studies in Linguistics* 38, 471-492.
- Russian National Corpus. (2016). Retrieved from <http://www.ruscorpora.ru/>
- Vasilieva, E. A. (2011). State power mythology: A sociological aspect. *Science and Education*, (3), 119-121.
- Van Valin, R.D. Jr. & R. J. LaPolla, (1997). *Syntax, Structure, Meaning and Function*. Cambridge: Cambridge University Press.