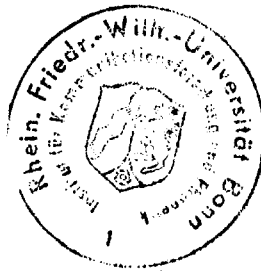


SUBCLASSIFICATION OF PARTS OF SPEECH IN RUSSIAN: VERBS*

A. Andreyewsky

International Business Machines Corporation
Thomas J. Watson Research Center
P. O. Box 218
Yorktown Heights, New York, 10598



*This work was partly sponsored by the Information Processing Laboratory, Rome Air Development Center, United States Air Force, under Contract AF 30(602)-3301

Abstract

In a trial study, about 500 Russian verbs were coded using 44 potential classificatory criteria. Through sorting and the introduction of a metric, numerous groupings were obtained. Initial results suggest that, with proper refinements, the approach described can provide useful information that may be employed in syntactic analysis and certain information retrieval applications.

0.0 Introduction

As part of a broader effort to extend the existing traditional part-of-speech classification in modern Russian, this study of verbs is oriented toward developing an improved basis for syntactic analysis. Moreover, it is hoped that the refinements introduced will be of interest in content analysis. To this end, an extensive set of potential classificatory criteria has been selected, in the hope that eventually this categorization can be optimized and extended to other parts of speech.

1.0 The Experiment

The 514 verbs analyzed came from two sources: (a) a randomized sample of 370 entries (1) and (b) a list of the most frequently used Russian verbs (2) from which the first 144 entries were selected.

The classificatory criteria, subdivided into two groups, are discussed in Section 1.1 below. Generally, each verb was taken in a particular meaning (stirat', for instance, as "to erase" and not as "to launder") and English equivalents used solely for purposes of identification. At the same time, for reasons of convenience, provisions were made in coding to allow for coexisting alternatives. Thus, for properties A and B there can be four possibilities which are represented by the following numerical codes: 1 - "A", 2 - "B", 3 - "AB", 0 - "neither applies".

After the verbs and appropriate codes were punched on cards, verbs with identical codes were compared. To obtain additional clustering, a program, written by R. F. Hubbard for the IBM 7040, compared the code vector of each card against those of the rest of the sample. The distance between any two entries was calculated by taking the square root of the sum of the squares of distance between corresponding positions in their code vectors as defined by the following table:

$$\begin{array}{l}
 0 \rightarrow 0 = 0 \quad 0 \rightarrow 2 = 4 \quad 1 \rightarrow 1 = 0 \quad 1 \rightarrow 3 = 1 \quad 2 \rightarrow 3 = 1 \\
 0 \rightarrow 1 = 4 \quad 0 \rightarrow 3 = 6 \quad 1 \rightarrow 2 = 2 \quad 2 \rightarrow 2 = 0 \quad 3 \rightarrow 3 = 0
 \end{array}$$

1.1 Tests

Since one of the main objectives of this study has been to establish the relevance of various classificatory criteria, these were tested in two groups as described below. The selection of criteria, based on studies of existing grammars of Russian, was directed toward discovering solutions for problems arising or likely to arise in machine-assisted syntactic analysis.

1.1.1 Test I

In this test, the verbs were coded according to their ability to combine with selected prepositional phrases, certain adverbs, and the chto-introduced object clauses. Most of the examples are derived from the discussion of slovosochetaniye (grammatically bound word group) problem in the Academy Grammar (3). While the English meanings supplied do reflect certain semantic differences the main objective has been to test not only the ability of a given verb to co-occur with certain types of phrases (examples are used solely for illustration) or classes of adverbs but to trace what effect the verb has on their syntactic function, if any.

1.1.1.1 Classificatory Criteria

- | | | |
|---|---|---|
| 1) ... <u>do menya</u>
(A) before me
(B) as far as me | 4) ... <u>k mitingu</u>
(A) for the ...
(B) to the ...
meeting | 7) ... <u>u Ziny</u>
(A) at Zina's
(B) from Zina |
| 2) ... <u>do rassveta</u>
(A) before dawn
(B) until dawn | 5) ... <u>k nam</u>
(A) to us
(B) toward us | 8) ... <u>pod kapustu</u>
(A) for cabbage
(B) under cabbage |
| 3) ... <u>iz-za stola</u>
(A) because of ...
(B) from behind
the table | 6) ... <u>za obedom</u>
(A) after (to get)...
(B) during dinner | 9) ... <u>za stol</u>
(A) at the
table
(B) behind the
table |

- | | | |
|---|--|--|
| <p>10) ... <u>za brata</u>
(A) in brother's place
(B) for brother's sake</p> <p>11) ... <u>po oshibke</u>
(A) a mistake apiece
(B) by mistake</p> | <p>12) ... <u>yashchik iz-pod uglya</u>
(A) coal crate
(B) crate from under the coal</p> <p>13) ... <u>o stol*</u>
against the table</p> <p>14) ... <u>po vodu*</u>
to get water</p> | <p>15) ... <u>chto napishet*</u>
that + (subject) + will write</p> <p>16) ... <u>nadvoe*</u>
in two (as in cutting)</p> <p>17) ... <u>ochen'*</u>
very much</p> <p>18) ... <u>o sestre*</u>
about the sister</p> |
|---|--|--|

1. 1. 1. 2 Results of Sorting

Sorting revealed some of the following groupings with identical codes:

- | | | |
|---|---|---|
| <p>A1 <u>zakipet'</u>
(to boil)
<u>prosolit'sya</u>
(to turn salty)</p> <p>A2 <u>vzdrognut'</u>
(to flinch)
<u>ustavat'</u>
(to become tired)
<u>ustat'</u>
(to become tired)
<u>izzyabnut'</u>
(to become chilled)</p> <p>A3 <u>sovrat'</u>
(to tell a lie)
<u>soobrazit'</u>
(to grasp)
<u>dogadat'sya</u>
(to surmise)</p> | <p><u>nameknut'</u>
(to hint)</p> <p>A4 <u>raznuzdat'</u>
(to let become undisciplined)
<u>vospitat'</u>
(to educate)</p> <p>A5 <u>vychest'</u>
(to subtract)
<u>izderzhat'</u>
(to spend)</p> <p>A6 <u>otrubit'</u>
(to chop off)
<u>vskryt'</u>
(to open up)</p> <p>A7 <u>nabryzgat'</u>
(to sprinkle on)
<u>rasprostranit'</u>
(to spread)</p> | <p>A8 <u>vbezhat'</u>
(to run in)
<u>yavit'sya</u>
(to appear)</p> <p>A9 <u>podumat'</u>
(to think)
<u>bredit'</u>
(to rave)
<u>okhat'</u>
(to moan)</p> <p>A10 <u>gordit'sya</u>
(to be proud)
<u>veselit'sya</u>
(to enjoy self)
<u>voskhishchatsya</u>
(to admire)</p> <p>All <u>verit'</u>
(to believe)
<u>toskovat'</u>
(to be sad, to pine)</p> |
|---|---|---|

* Only test ability to combine in the meaning indicated.

	<u>grustit'</u> (to be sad, to yearn)	A16 <u>nastroit'</u> (to incite)	A21 <u>morosit'</u> (to drizzle)
	<u>skuchat'</u> (to be bored)	<u>bespokoit'</u> (to disturb)	<u>nakrapivat'</u> (to sprinkle)
	<u>fantazirovat'</u> (to dream)	<u>obizhat'</u> (to offend)	<u>poroshit'</u> (to snow)
A12	<u>volnovat'sya</u> (to worry)	<u>proklinat'</u> (to damn)	A22 <u>farshirovat'</u> (to stuff)
	<u>opasat'sya</u> (to be afraid)	<u>portit'</u> (to spoil, ruin)	<u>sintezirovat'</u> (to synthesize)
A13	<u>zapryatyvat'</u> (to hide)	A17 <u>bakhvalit'sya</u> (to brag)	A23 <u>klassifitsirovat'</u> (to classify)
	<u>vovlekat'</u> (to draw in)	<u>likovat'</u> (to rejoice)	<u>razbivat'</u> (to break)
A14	<u>berech'</u> (to save)	A18 <u>razrykhlyat'</u> (to loosen)	A24 <u>begat'</u> (to run)
	<u>poberech'</u> (to save)	<u>razdrobit'</u> (to pulverize)	<u>prikhodit'</u> (to come)
	<u>uderzhat'</u> (to withhold)	A19 <u>besedovat'</u> (to converse)	A25 <u>nesti</u> (to carry)
A15	<u>vzgromozdit'sya</u> (to perch self)	<u>soveshchat'sya</u> (to confer)	<u>vezti</u> (to cart)
	<u>uglubit'sya</u> (to go deep into)	A20 <u>razodrat'</u> (to tear)	A26 <u>volochit'</u> (to drag)
	<u>rassazhivat'</u> (to seat)	<u>rasshibit'</u> (to break, bust)	<u>tashchit'</u> (to pull)
			A27 <u>doyti</u> (to reach, (walking))
			<u>doletet'</u> (to reach (flying))

1. 1. 1. 3 Results of the Introduction of the Metric

On the basis of preliminary results, the maximum distance considered was set at 10. Given this arbitrary limitation, the metric produced various groupings. The majority of them contained some "noise" - i. e. , apparently incorrect entries were brought together or several distance groupings turned out insufficiently differentiated. Partly responsible for this are: the method employed, the distances selected, and the occasional errors that crept in during the analysis and subsequent processing. These factors are discussed in greater detail below (1. 1. 1. 4).

Some of the more interesting outcomes were as follows:

- A28 Groups A11 (verit', toskovat', grustit', skuchat', and fantazirovat'), A17 (bakhvalit'sya and likovat'), A12 (volnovat'sya and opasat'sya), and the verb bespokoit'sya (to worry)
- A29 Group A16 (nastroit', bespokoit', obizhat', proklinat', portit') and the verb nenavidet' (to hate).
- A30 Group A10 (voskhishchiat'sya, veselit'sya, gordit'sya) and verbs vozmudit'sya (to become disgusted) and boyat'sya (to be afraid).
- A31. Group A8 (yavit'sya, vbezhat'), the following verbs: vernut'sya (to return), prikhodit' (A24), begat' (A24), vyyti (to step out), podyezhat' (to drive up), yezdit' (to ride), vyekhat' (to go away), kinut'sya (to lunge), vypolzti (to crawl out), doletet' (A27), and doyti (A27).
- A32 garantirovat' (to guarantee), pokazyvat' (to show), demonstrirovat' (to demonstrate)
- A33 sovrat' (to lie), poverit' (to believe), uverit' (to assure)
- A34 znat' (to know), ozhidat' (to expect), videt' (to see).
- A35 nagryanut' (to come unexpectedly), zaekhat' (to stop by), probezhat'sya (to run), otstupit' (to retreat).

1.1.1.4 Comments

The problems stemming from the application of the metric (the "numbers game") mentioned in 1.1.1.3 reflect a characteristic of statistical inference jocularly compared by an anonymous author to a bikini bathing suit: being sufficiently suggestive, but not revealing. In this regard, alternative approaches have been considered and will be tried in the near future. As it turned out in practice, however, the metric did provide useful insights which can point the way toward developing a more powerful set of classificatory criteria. This, in turn, can foster increased reliance on simple sorting procedures based on proper ranking and grouping of the criteria themselves.

While not unexpectedly, the verbs of motion in the broad sense of the term came out more clearly in the classification than did any other groups, interesting subclasses of abstract verbs, exhibiting unexpected shades of valuation also emerged.

1. 1. 2 Test II

In contrast to Test I, this test placed a relatively lesser emphasis on syntagmatic relationships and stressed a mixture of formal and semantic properties. On the whole, except where noted, the two tests were developed independently of one another. While Test I was based on materials derived from the Academy Grammar of Russian (3), Test II benefited from experience gained in dealing with the problems encountered in machine translation output and from studies conducted preparatory to launching syntactic analysis.

1. 1. 2. 1 Classificatory Criteria

In view of the extensive nature of this test, the description of various criteria used is given here in abbreviated notation.

- | | | |
|---|--|---|
| 1) (A) imperfective
(B) perfective | 7) passive participle:
(A) past
(B) present | 14) meaning affected by
(A) governed infinitive
(B) object(s) |
| 2) verb (1/3) or
"verboïd" (2/0);
"concrete" (1/0)
or "abstract"
(2/3): when "yes"
answer is possible
under 1. 1. 1. 1. 17. | 8) gerundial forms:
(A) present
(B) past | 15) subject preference:
(A) inanimate
(B) animate |
| 3) <u>is given</u> form
(A) reflexive
(B) non-reflexive | 9) action (gerund):
(A) parallel
(B) sequential | 16) verb governs:
(A) infinitives
(B) objects |
| 4) <u>generally</u> :
(A) non-reflexive
(B) reflexive | 10) deverbal nouns:
(A) in <u>-enie, -ka</u>
(B) other forms | 17) object preference:
(A) animate
(B) inanimate |
| 5) when reflexive,
meaning:
(A) active
(B) passive | 11) deverbal nouns:
(A) concrete
(B) abstract | 18) (A) motion verb
(broad sense)
(B) action perceived |
| 6) participial forms:
(A) active
(B) passive | 12) verb used:
(A) personally
(B) impersonally | 19) verb describes:
(A) action
(B) state |
| | 13) verb function:
(A) link, auxiliary
(B) other | 20) (A) beginning
(B) end of action |

- | | | |
|--|--|--|
| 21) verb is one of:
(A) being
(B) becoming | 23) action directed:
(A) downward
(away)
(B) upward
(toward) | 25) reference to:
(A) duration
(B) intensity |
| 22. action described:
(A) outward-
(B) inward-
directed | 24) action in respect
to object:
(A) contacts
(B) permeates | 26) action produces:
(A) decrease
(B) increase |
| | | 27) action describes:
(A) gain
(B) loss |

1. 1. 2. 2 Results of Sorting

The following groupings had identical codes:

- | | | |
|--|--|---|
| B1 <u>skuchat'</u> (A11)
(to be bored)
<u>toskovat'</u> (A11)
(to be sad) | B4 <u>izderzhat'</u>
(to expend)
<u>istratit'</u>
(to spend) | <u>navyazyvat'</u>
(to tie on)
<u>skladyvat'</u>
(to put together) |
| B2 <u>morosit'</u> (A21)
(to drizzle)
<u>poroshit'</u> (A21)
(to snow) | B5 <u>nosit'</u>
(to carry)
<u>tashchit'</u>
(to pull)
<u>volochit'</u>
(to drag) | B7 <u>ozhivit'</u>
(to vivify)
<u>uverit'</u>
(to assure) |
| B3 <u>nakrapyvati'</u> (A21)
(to sprinkle)
<u>mertsat'</u>
(to twinkle) | B6 <u>podshivat'</u>
(to attach) | |

1. 1. 2. 3 Results of the Introduction of the Metric

Comments made in 1. 1. 1. 3 above, apply. Because of a greater number of classificatory criteria the results of introducing the metric were more important in this test. Numbers in parentheses preceding each verb indicate distances from the first verb in the group.

- | | | |
|---------------------------------------|-------------------------------------|------------------------------------|
| B8 <u>pridavit'</u>
(to squeeze) | B9 <u>vosstat'</u>
(to riot) | B10 <u>prichesat'</u>
(to comb) |
| (1) <u>prishchemit'</u>
(to pinch) | (1) <u>vystupit'</u>
(to appear) | (1) <u>zaputat'</u>
(to tangle) |

B11 <u>vbezhat'</u> (to run in)	B17 <u>otkryt'</u> (to open)	B23 <u>nestis'</u> (to dash)
(1) <u>vypolzti</u> (to crawl out)	(5) <u>ubavit'</u> (to decrease)	(7) <u>bezhat'</u> (to run)
B12 <u>napevat'</u> (to hum)	B18 <u>vynesti</u> (to carry out)	B24 <u>vydelit'</u> (to single out)
(2) <u>veshchat'</u> (to speak with authority)	(5) <u>vypustit'</u> (to let out)	(7) <u>vypisat'</u> (to write out)
B13 <u>temnet'</u> (to grow dark)	B19 <u>zheltet'</u> (to turn yellow)	B25 <u>potusknet'</u> (to dull)
(2) <u>teplet'</u> (to grow warm)	(5) <u>umirat'</u> (to die)	(7) <u>zatverdet'</u> (to harden)
B14 <u>vyrabotat'</u> (to develop)	B20 <u>terrorizirovat'</u> (to terrorize)	B26 <u>prikrepit'</u> (to fasten)
(3) <u>vyuchit'</u> (to learn)	(5) <u>khvalit'</u> (to praise)	(8) <u>nav'yuchit'</u> (to pack on)
B15 <u>khmurit'sya</u> (to frown)	B21 <u>viset'</u> (to hang)	B27 <u>vozvratit'</u> (to return)
(3) <u>tumanit'sya</u> (to grow gloomy)	(6) <u>lezhat'</u> (to lie)	(8) <u>dopolnit'</u> (to augment)
B16 <u>razbushevat'sya</u> (to start raging)	B22 <u>podognat'</u> (to drive up)	B28 <u>vvesti</u> (to introduce)
(4) <u>uchastit'sya</u> (to become more frequent)	(6) <u>navestit'</u> (to visit)	(9) <u>dobavit'</u> (to add)

In addition to shorter groups described above, longer groupings were observed. Thus, otdokhnut' (to rest) (8) utikhnut' (to quiet down), and (10) ugasnut' (to become extinguished) or nabryzgat' (to sprinkle on), (2) nakinut' (to throw on), (3) vzvalit' (to pile on), and (4) nastrocit' (to sew on) are some of the examples.

In other cases, apparently incongruous groups like the following: strekotat' (to chirr), (1) moshennicat' (to swindle), (5) fokusnichat' (to juggle), (5) nakrapivat' (to sprinkle), (5) mertsat' (to twinkle) (6) zvenet' (to ring) emerged. However, upon closer examination it became apparent that nakrapivat', mertsat', and zvenet' fall in a group clearly distinguishable from the one containing the other verbs. Further, fokusnichat' and zvenet' showed sufficient distance within respective groups suggesting at least four different basic groups in all.

1. 1. 2. 4 Comments

Aside from the problems traceable to statistics, the sets of criteria selected for Test II are more open to debate than those found in Test I. However, correlations between both tests indicate that some of the criteria are relevant and that others are, at least, redundant. As observed from minor differences in two versions of coding of nine verbs introduced six months apart, the results of Test II are less reliable.

1. 1. 3 Comparison of Test I and Test II

As noted in 1. 1. 2 above, the two tests differ in the base from which they were derived. Accordingly, the results obtaining from Test I are both intuitively and actually more reliable. Yet, as suggested in 1. 1. 1. 4, to the extent that the results of the application of the metric tend to supplement sorting, the results of Test II tend to back up many of the findings of Test I.

Given a small sample, it is difficult to make any generalizations. At the same time, the evidence emerging so far suggests some subtle differences in the two tests. Basically, in both cases the results of the metric application show little or no discrimination between antonyms. However, the groupings resulting from Test II tend to be, if at all, held together by similarity of content. The results of Test I, in contrast, have a peculiar sort of outward, formal similarity in the manifestation of processes described by the verbs in question.

2. 0 The Outlook

In the months ahead, it is hoped that the small corpus can be increased and the time required to code each entry reduced to reasonable proportions. While in many respects the results of both tests are self-proving, rigorous evaluation criteria will have to be formulated in detail.

As far as potential application of the results obtained is concerned, especially the information derivable from Test I could be immediately put to use to improve (together with classification of nouns currently in progress) the translation of verb-governed prepositional phrases. It is likely that this syntagmatic patterning will extend to larger structures dominated by the verb. Further, if the apparent trends persist, some framework of semantic classification can be anticipated. To what extent this will be possible to accomplish by computers alone and the degree to which such a classification

will satisfy the needs of computer processing remains to be established. While it can be argued that any classification is likely to produce some classes, we take solace in the fact that the methodology employed even in such classics as Roget's Thesaurus remains unknown to this day.

Sources

1. This sample was selected from the Daum and Schenck Dictionary in another connection and was generally random in its intent more than its methodology.
2. A. K. Demidova, O. G. Motovilova, G. D. Shevchenko, E. P. Chaplygin, Naiboleye upotrebitel'nyye glagoly sovremennogo russkogo yazyka (The Most Frequently Used Verbs in Modern Russian), Moscow, USSR Academy of Sciences Publishing House, 1963.
3. V. V. Vinogradov, ed., Grammatika russkogo yazyra (Grammar of the Russian Language) Moscow, USSR Academy of Sciences Publishing House, 1960, Vol. II, Part I, pp. 113-230.