

1 Art. No.: 831

2

3 **Making choices in Russian: Pros and cons of statistical**
4 **methods for rival forms**

5 Выбор вариантных форм в русском языке: плюсы и минусы
6 различных моделей статистического анализа

7

8 Making choices in Russian

9

10 R. Harald Baayen

11 University of Tübingen, Tübingen, Germany

12 E-Mail: harald.baayen@uni-tuebingen.de

13

14 Anna Endresen

15 University of Tromsø, Tromsø, Norway

16 E-Mail: anna.endresen@uit.no

17

18 Laura A. Janda

19 University of Tromsø, Tromsø, Norway

20 E-Mail: laura.janda@uit.no

21

22 Anastasia Makarova

23 University of Tromsø, Tromsø, Norway

24 E-Mail: anastasia.makarova@uit.no

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Tore Nessel
University of Tromsø, Tromsø, Norway
E-Mail: tore.nessel@uit.no

Corresponding author: Laura Janda

Abstract Sometimes languages present speakers with choices among rival forms, such as the Russian forms *ostrič* vs. *obstrič* ‘cut hair’ and *proniknuv* vs. *pronikši* ‘having penetrated’. The choice of a given form is often influenced by various considerations involving the meaning and the environment (syntax, morphology, phonology). Understanding the behavior of rival forms is crucial to understanding the form-meaning relationship of language, yet this topic has not received as much attention as it deserves. Given the variety of factors that can influence the choice of rival forms, it is necessary to use statistical models in order to accurately discover which factors are significant and to what extent. The traditional model for this kind of data is logistical regression, but recently two new models, called ‘tree & forest’ and ‘naive discriminative learning’ have emerged as alternatives. We compare the performance of logistical regression against the two new models on the basis of four datasets reflecting rival forms in Russian. We find that the three models generally provide converging analyses, with complementary advantages. After identifying the significant factors for each dataset, we show that different sets of rival forms occupy different regions in a space defined by variance in meaning and environment.

50 **Аннотация** Носители языка часто сталкиваются с ситуацией выбора
51 вариантных форм, таких как рус. *остричь* и *обстричь* или *проникнув* и
52 *проникши*. На выбор варианта могут влиять различные факторы, включая
53 семантику и контекстное окружение (синтаксическое, морфологическое и
54 фонологическое). Изучение поведения вариантных форм необходимо для
55 понимания соотношения означающего и означаемого в языке, однако этот
56 вопрос до сих пор не получил должного внимания. Ввиду того, что выбор
57 вариантной формы может зависеть от факторов различного рода,
58 необходимо использовать методы статистического анализа: они
59 позволяют точно определить, какие факторы являются главными и какова
60 доля их влияния. Обычно для такого типа языковых данных применяется
61 модель логистической регрессии, однако недавно появились две
62 альтернативные модели—‘случайный лес’ и ‘наивное различительное
63 обучение’. Мы сравнили эффективность логистической регрессии и двух
64 новых моделей статистического анализа на материале четырех баз
65 данных, собранных для ряда вариантных форм русского языка. Все три
66 модели дают в целом схожие результаты, но каждая имеет свои
67 преимущества. В статье выявлены определяющие факторы для каждого
68 набора данных, а также показано, что исследованные нами вариантные
69 формы размещаются в различных зонах системы двух осей координат—
70 оси различия по значению и оси различия по контекстным условиям.

71

72 Keywords: rival forms, statistical models, form-meaning relationship

73

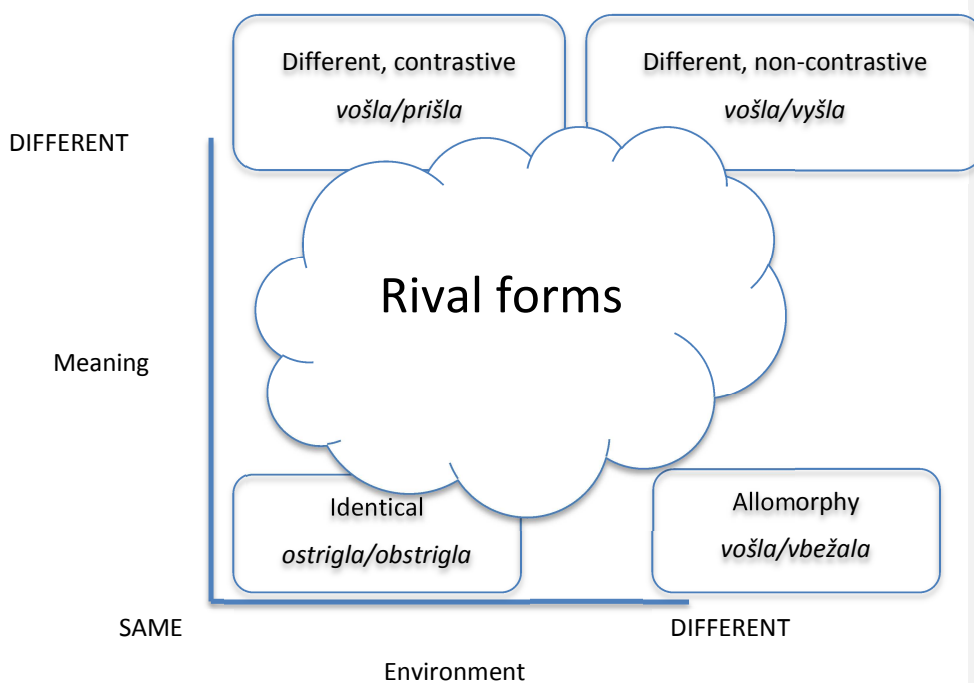
74 **1 Introduction**

75 This article focuses on the statistical analysis of rival forms in language. Rival
76 forms exist when a language has two (or more) forms that express a similar
77 meaning in similar environments, giving the speaker a choice of options. The
78 choice made between rival forms is often influenced by a range of factors such
79 as the syntactic, morphological, and phonological environment. We will
80 commence by examining the place of rival forms in the form-meaning
81 relationship.

82 The form-meaning relationship is essential to language, yet highly complex,
83 both in terms of the relationship itself, and in terms of the environments in
84 which this relationship is valid. We can think of this relationship as a three-
85 dimensional space, with form, meaning, and environment as the three axes that
86 define this space. Each axis has a continuum of values that range from perfect
87 identity (when the form, meaning, and environment are exactly the same) to
88 contrast (when the form, meaning, and environment are entirely different). At
89 these two extremes we have trivial cases of either identical items (with
90 identical meanings found in identical environments), or different items (with
91 different meanings found in different environments). However, each axis
92 captures a gradient that also includes variants that lie between identity and
93 difference, involving near-identity, similarity, overlap, and varying degrees of
94 contrast, fading to mere (non-contrastive) difference. If we choose to look only
95 at cases showing difference in form, then meaning and environment yield a
96 two-dimensional space, as visualized in Fig. 1.

97 In addition to the labels at the four corners of Fig. 1, synonymy lies along
98 the bottom horizontal axis of the space. Whereas strictly speaking synonyms

99 should have the ‘same’ meaning, in reality even the best of synonyms are
 100 usually near-synonyms, with slightly different shades of meaning. Thus
 101 synonymy is a gradient phenomenon, with some synonyms overlapping nearly
 102 entirely in terms of both meaning and environment, but others showing some
 103 deviation.¹ The space in the center of Fig. 1 is labeled ‘Rival forms’ and
 104 includes relationships involving near-synonymy and partial synonymy as well
 105 as various degrees of overlap in terms of environments.



116 **Fig. 1 The space defined by variance in meaning and environment**

117
 118 Linguists tend to focus on the four corners of this space, which we can
 119 illustrate with Russian verbal prefixes and environments involving syntactic,
 120 morphological (word-formation), and phonological factors. Let’s begin at the
 121 origin, where the environment and meaning are the same, and then continue in

¹ This article does not address antonyms, which are actually very similar to synonyms, providing contrast in only one (or a few) parameters, but usually found in the same environments and thus located along the leftmost vertical axis of Fig. 1.

122 a clockwise direction around the corners from there. For example, if we have
123 two attestations *mat' ostrigla volosy rebenku* and *mat' obstrigla volosy rebenku*
124 'the mother cut the child's hair', we have the same meaning and the same
125 environment (in terms of word-formation and syntax), and the variant forms
126 *o-* and *ob-* perform identical roles; for this example the prefixes are in free
127 variation. If we change the meaning, but keep the same word-formation and
128 syntactic environment, we then get contrasting meanings of the prefixes
129 *vo-* and *pri-* as in *mat' vošla v cerkov'* 'mother entered (into) the church' and
130 *mat' prišla v cerkov'* 'mother came to church', where the former phrase
131 emphasizes the church as a building and the latter one refers to a functional
132 relationship (it is most likely that mother in this phrase is attending a service or
133 other meeting). The fact that *vo-* and *pri-* can occur in some of the same
134 environments makes it possible for their meanings to be used contrastively.
135 Next is a case where both the meaning and the environment (in terms of
136 syntax) are different, as in *mat' prišla v cerkov'* 'mother entered (into) the
137 church' and *mat' vyšla iz cerkvi* 'mother exited (from) the church', where the
138 prefixes *vo-* and *vy-* are simply different in both their meaning and their
139 distribution. In the last corner we find allomorphy, traditionally defined as a
140 relationship of different forms that share a meaning but appear in
141 complementary distribution (Bauer 2003, p. 14; Booij 2005, p. 172;
142 Haspelmath 2002, p. 27; Matthews 1974, p. 116). Here we have phonologically
143 conditioned examples like *mat' vošla v cerkov'* 'mother entered (into) the
144 church (walking)' and *mat' vbežala v cerkov'* 'mother entered (into) the church
145 (running)', where *vo-* and *v-* are allomorphs and their different distribution is
146 conditioned by the phonological shape of the root to which they are attached.

147 Here the environment is phonological instead of being an environment that
148 involves word-formation or syntax.

149 The space between the four points in Fig. 1 has not been thoroughly
150 explored by linguists, yet arguably contains many of the most interesting form-
151 meaning-environment relationships found in language. Although rival forms
152 have received some attention in the literature (cf. Aronoff 1976 and Riddle
153 1985 on rival affixes in English word-formation, such as *-ity* and *-ness*), this is
154 an understudied topic. More empirical studies are needed. The present article is
155 an attempt to fill this need.

156 We examine four cases, all of which involve a choice between two rival
157 forms: 1) *gruzit'* 'load' and its prefixed perfective forms which appear in two
158 rival constructions, 2) the prefixes *pere-* vs. *pre-*, 3) the prefixes *o-* vs. *ob-*, and
159 4) the use of *nu-* vs. \emptyset -forms of verbs like *(ob)soxnut'* 'dry'. Although this is
160 primarily a methodological article, the case studies all relate to the topic of this
161 special issue, namely the understanding of time in Russian since they involve
162 rival forms of Russian verbs associated with perfectivizing prefixes and the *nu*
163 suffix. Each case study is supported by an extensive dataset and a variety of
164 statistical models are applied in order to discover the complex structures in the
165 form-meaning-environment relationships. Section 2 provides a general
166 discussion of the range of options for statistical analysis and problems posed by
167 various datasets. The studies are presented in Sect. 3, which relates each case
168 study to the parameters in Fig. 1 and also states the linguistic objective of each
169 study. The results of the analyses are summarized in the conclusions in Sect. 4.
170 All the datasets and the code used for their analyses are available at this site:
171 ansatte.uit.no/laura.janda/RF/RF.html. All analyses are performed using the

172 statistical software package R (2011), which is available for free at [www.r-](http://www.r-project.org)
173 [project.org](http://www.r-project.org).

174

175 **2 Options for statistical analysis**

176 This section presents the three statistical models that we have compared: the
177 logistic regression model, the tree & forest model (combining classification
178 trees with random forests), and the naive discriminative learning (NDL) model.

179 Despite the variety of data represented in our four case studies, they share a
180 similar issue: each one presents a pair of rival forms and their distribution with
181 respect to an array of possible predicting factors. If we call the rival forms X
182 vs. Y , then we can define a categorical factor, say Prefix, that has as its levels
183 two rival forms, the prefixes X and Y . Given semantic and environmental
184 predictors such as Aspect, Animacy, Frequency, etc., we can restate all of the
185 case studies in terms of questions like these:

186

- 187 1. Which combinations of values for Aspect, Animacy, Frequency, etc.,
188 predict the value of the response variable Prefix?
- 189 2. How do the predictors rank in terms of their relative strength or
190 importance?
- 191 3. If we build a model that optimizes the use of the predictors to predict
192 the response (X vs. Y), how accurate is that model, how well does it
193 capture valuable generalizations without being overly affected by low
194 level variation that is merely noise?

195

196 We can think of these questions as being parallel to many other types of
197 questions one might ask in many non-linguistic situations such as:

198

- 199 • Predicting whether patients will get cancer ($X = \text{yes}$ vs. $Y = \text{no}$) given
200 possible predictors such as age, body mass index, family history,
201 smoking history, alcohol use, diet, exercise, etc.
- 202 • Predicting which candidate voters will select ($X = \text{democrat}$ vs. $Y =$
203 republican) given possible predictors such as age, race, religion,
204 income, education level, region, etc.
- 205 • Predicting which product ($X = \text{name brand}$ vs. $Y = \text{generic brand}$)
206 consumers will select given possible predictors such as price, volume,
207 advertising, packaging, etc.

208

209 The popular method statisticians apply to such situations with a binary
210 response variable is logistic regression (cf. Baayen 2008, Chapter 6). The first
211 subdiscipline in linguistics to make use of logistic models is sociolinguistics
212 (Cedergren and Sankoff 1974; see also Tagliamonte and Baayen 2012). More
213 recently, this type of modeling has also been applied to lexical choices (Arppe
214 2008) and grammatical constructions (Bresnan, Cueni, Nikitina and Baayen
215 2007). The strategy of a regression model is to model the functional
216 relationship between the response and its predictors as a weighted sum
217 quantifying the consequences of changing the values of the predictors. For
218 factorial predictors (such as perfective versus imperfective), the model
219 specifies the change in the group means when going from one factor level (e.g.
220 perfective) to the other (imperfective). For numerical predictors (such as

221 frequency), the model specifies the consequences of increasing the predictor's
222 value by one unit. The goal of a logistic regression model is to predict the
223 probability that a given response value (X , or alternatively, Y) will be used. It
224 does so indirectly, for mathematical reasons, by means of the logarithm of the
225 odds ratio of X and Y . The odds ratio is the quotient of the number of
226 observations supporting X and the number of observations supporting Y . The
227 log of the odds ratio is negative when the count for Y is greater than the count
228 for X . It is zero when the counts are equal. It is positive when the counts for X
229 exceed the counts for Y .

230 Fitting a logistic regression model to the data amounts to finding the
231 simplest yet most adequate model for the data. A model is simpler when it has
232 fewer predictors. A model is more adequate when its predictions approximate
233 the observations more closely. Typically, one will have to find a balance
234 between the two, by removing predictors that do not increase the goodness of
235 fit, and by adding in predictors that make the model more precise. In the
236 present study, we use a hypothesis-driven search for the best model.

237 An important concept in statistical modeling is that of an interaction
238 between predictors. Consider two predictors, for instance, Animacy (with
239 levels animate and inanimate) and Aspect (with levels perfective and
240 imperfective). There is no interaction when a change in Animacy (or a change
241 in Aspect) is the same for all the levels of the other factor. However, when the
242 likelihood of response X increases when changing from animate to inanimate
243 for perfective verbs, but decreases (or increases less) for imperfective verbs,
244 then an interaction of Animacy by Aspect is at issue. Adding in interaction
245 terms may substantially increase the goodness of fit of a model.

246 The output of a logistic regression model gives us information that
247 addresses all three questions stated above:

248

249 1. We can discover which of the predictors predict the value of the
250 response variable by checking whether a change in the value of a given
251 predictor implies a significant change in the value of the response. In
252 the case of logistic regression, this implies a significant change in the
253 value of the log-odds, which translates into a significant change in the
254 probability of, e.g., the response value X .

255 2. Information about the relative strength and importance of a predictor
256 can be obtained by inspecting both the magnitude of its effect on the
257 response, and by considering the extent to which adding the predictor to
258 the model increases its goodness of fit. This is typically accomplished
259 with the Akaike's Information Criterion (AIC) measure. Lower values
260 of AIC indicate a better model fit.

261 3. It is possible to evaluate the accuracy of the model by comparing its
262 predictions (whether the response has as its value X or Y) with the
263 actual observed values. Accuracy measures can be imprecise, however,
264 because the model delivers probabilities whereas the observations are
265 categorical (X or Y). One can posit that a probability of X greater than
266 or equal to 0.5 is an X response, and a probability of X less than 0.5 a Y
267 response. But this procedure makes it impossible to see that the model
268 might be correctly predicting differences in probability below (or
269 above) 0.5. For instance, changing from inanimate to animate might
270 raise the probability of an X response from 0.6 to 0.8. The accuracy

271 measure cannot inform us about this. A second measure, C , the index of
272 concordance, has been developed that does not have this defect, and
273 therefore provides a more precise measure of how well the model
274 performs. For a model to be considered a good classifier, the value of C
275 should be at least 0.8.

276

277 Most readers who are not already proficient with statistics are likely to express
278 frustration at this point, since the tasks of designing an optimal logistic
279 regression model and then interpreting the output are rather daunting. In fact,
280 guidelines and principles for finding the optimal model are an active area of
281 research, with computer scientists proposing algorithms that will find the best
282 fitting model on the one hand, and researchers preferring hypothesis-driven
283 model selection on the other hand. The goal of this article is to illustrate
284 logistic modeling, but to complement it with two alternative models that are
285 more straightforward to use, and that sometimes yield results that are more
286 intuitive in their interpretation. The two alternatives we present here are: 1.
287 classification trees and random forests (tree & forest, cf. Strobl, Malley and
288 Tutz 2009) and 2. NDL (Baayen 2011). Both alternatives eliminate the step of
289 searching for an optimal regression model: They arrive at their optimal
290 solutions on their own. Especially in the case of the tree & forest method, the
291 output is often easier to interpret as well: The classification tree is an entirely
292 intuitive diagram of the outcomes that are predicted and yielded by various
293 combinations of predictor values.

294 Logistic regression modeling is a very powerful tool when the data do not
295 violate certain underlying mathematical assumptions. One such assumption is

296 that when testing for interactions between two factors, all combinations of
297 factor levels should be attested. For linguistic datasets, this condition is not
298 always satisfied, often because the grammar does not allow for certain
299 combinations. For instance, in the *nu* vs. \emptyset dataset, there are no unprefix past
300 gerunds. An advantage of classification trees & random forests and NDL is that
301 they do not impose distributional constraints, and are thus better suited for
302 many types of datasets involving naturalistic data on rival linguistic forms.

303 In the R programming environment, all three types of models use the same
304 basic format for the formula that relates the rival forms to the predictors. This
305 formula places the predicted variable to the left of a tilde \sim and places the
306 predictors to the right, separated by plus '+' signs.² Our abstract and
307 hypothetical examples above would be rendered by these formulas (using
308 Response to refer to X vs. Y):

309

310 1. Rival linguistic forms:

311 Response \sim Aspect + Animacy + Frequency

312

313 2. Cancer prediction:

314 Response \sim Age + BodyMassIndex + FamilyHistory +

315 + SmokingHistory + AlcoholUse + Diet + Exercise

316

317 3. Voter choice prediction:

318 Response \sim Age + Race + Religion + Income +

319 + EducationLevel + Region

² The plus sign should be read as 'and' and not as a summation. It is only in the case of logistic models that the plus sign can be interpreted as summation, but then it indicates that the response is modelled as a weighted sum of the predictor values.

320

321 4. Consumer choice prediction:

322 Response ~ Price + Volume + Advertising +

323 + Packaging

324

325 While both the tree & forest model and NDL are non-parametric classification
326 models (as opposed to the parametric logistic model), they work on different
327 principles and this has implications for the kinds of datasets that can be
328 modeled and the results of analysis. The tree & forest model uses recursive
329 partitioning to yield a classification tree that provides an optimal partitioning of
330 the data, giving the best ‘sorting’ of observations separating the response
331 outcomes (see description of bootstrap samples below). It can literally be
332 understood as an optimal algorithm for predicting an outcome given the
333 predictor values.

334 NDL provides a quantitative model for how the brain makes the choice
335 between rival forms and constructions. This type of model makes use of a two-
336 layer network, the weights of which are estimated using the equilibrium
337 equations of Danks (2003) for the Rescorla-Wagner equations (Rescorla and
338 Wagner 1972) that summarize and bring together a wide body of results on
339 animal and human learning. The basic idea underlying this model is best
340 explained by an example. Consider English scrabble, and imagine a situation in
341 which one has a Q, an A, but no U. In that case, knowledge of the word *quid*,
342 an accepted word in English scrabble, will increase the chances of playing the
343 Q. The letter combination QA, although very infrequent, is an excellent cue for
344 the word *quid*. The greater the number of words with a given form pattern, the

345 less that form pattern will serve as a cue to the meaning of any specific word
346 with that pattern. NDL estimates from (corpus) data the strengths with which
347 form cues support a given meaning. Baayen, Milin, Đurđević, Hendrix &
348 Marelli (2011) showed that a simple naive discrimination network can account
349 for a wide range of empirical findings in the literature on lexical processing.
350 Baayen (2011) used a discrimination network to model the dative alternation in
351 English (Bresnan et al. 2007), and showed that such a network performed with
352 accuracy on a par with that of other well-established classifiers. This shows
353 that human probabilistic behavior can be understood as arising from very
354 simple learning principles in interaction with language experience as sampled
355 by corpus data. The NDL model can be pitted against naturalistic datasets in
356 order to ascertain to what extent human learning (under ideal conditions) and
357 statistical learning (using computational algorithms with no cognitive
358 plausibility) converge.

359 Both the tree & forest model and NDL provide a mechanism for validating
360 the model. These validation techniques assess how the results of a statistical
361 analysis will generalize to an independent dataset. Ideally one would build a
362 statistical model for a given phenomenon based on one dataset (the training
363 dataset) and then test the performance of that model using a second,
364 independent dataset (the validation dataset). In this way one can avoid circular
365 reasoning that would result from building and validating the model on the same
366 dataset (since of course the model will perform best if we ask it to predict the
367 outcomes of the data that were the input for its design). These techniques also
368 protect against overfitting the data. Overfitting occurs when the model reflects
369 variation that is characteristic of the particular sample of data, and this

370 interferes with how the model reflects the generalizations that are relevant to
371 the phenomenon under study in the population from which the data were
372 sampled. In other words, any given sample might misrepresent the relationship
373 between the rival outcomes and possible predictors due to chance variation,
374 and ideally this problem would be solved by using two samples, a training
375 dataset and an independent, new ‘validation’ dataset. Statisticians have
376 designed a variety of validation techniques in order to address the gap between
377 the ideal situation and the limitations of reality. In many cases it is not really
378 possible (or at least extremely difficult) to get two large independent samples
379 of the relevant data. Linguists face this problem, for example, due to limits on
380 corpus data: the size of any given corpus is finite, and once all the relevant data
381 from a given corpus has been mined out, it is not possible or very difficult to
382 get a second independent dataset that would be an equivalent sample in terms
383 of size and sources.

384 The basic idea underlying the validation techniques is to use part of the
385 available data for fitting (or training) the model, and the remaining part of the
386 data to test the predictions of the model on.

387 In the tree & forest model, bootstrap samples are used. A bootstrap sample
388 is a sample, drawn with replacement, of size N from a dataset with N
389 observations. As a consequence of replacement, some observations are
390 sampled more than once, and others are not sampled at all. The data points
391 sampled at least once constitute the in-bag observations on which we base
392 learning, the data points that are not sampled constitute the out-of-bag
393 observations, which we will predict.

394 NDL uses a ten-fold cross-validation. This validation technique partitions
395 the data into ten subsamples. Nine of the subsamples serve collectively as the
396 training dataset (the in-bag observations), while the remaining subsample is
397 used as a validation dataset (the out-of-bag observations on which we test our
398 predictions). This process is repeated ten times, so that each of the ten
399 subsamples has been used once as a validation dataset.

400 One thing to remember with both the random forest and NDL models is that
401 because randomization is used in the calculations, some of the output can differ
402 slightly each time these analyses are run. In fact, it is always a good idea to run
403 the validation procedure several times, to make sure that a particular result
404 does not depend on how the data happened to be sampled.

405 We will take up each dataset in turn, motivate our choice for the optimal
406 statistical model, and detail its interpretation. In addition to this primary goal of
407 alternative models and their interpretation, our secondary goal is to show how
408 statistical models can help us to explore and understand the structure of
409 naturalistic datasets such as the ones presented here. More specifically, we will
410 use statistical models as a sensitive multi-purpose tool for ferreting out the
411 relationships between rival forms and their predictors.

412

413 **3 Analyses**

414 The analyses are presented according to the relative complexity of the data,
415 starting with the most straightforward dataset. Each subsection below presents
416 a dataset by stating its name, source, overall size, rival forms, and values for
417 predictors. We then present the optimal statistical model and compare it with
418 other possible models and briefly discuss the results and what they tell us about

419 the rival forms and their behaviors. The first dataset is the one with the *gruzit'*
420 data (LOAD), which is relatively simple because it has few predictors, each
421 with few levels. This dataset is amenable to analysis by all three of the methods
422 we present in this article, yielding very similar results for all three. We give a
423 relatively detailed explanation of how to interpret the results of the three types
424 of models for the LOAD data and more abbreviated notes on the results for the
425 remaining datasets. Some additional details are available in the annotations to
426 the R script at ansatte.uit.no/laura.janda/RF/RF.html.

427

428 **3.1 *Gruzit'* and its perfectives in the theme – object vs. goal – object** 429 **constructions**

430 The objective of this case study is to show that so-called ‘empty’ perfectivizing
431 prefixes are actually distinct since they can show unique patterns of preference
432 for grammatical constructions. When prefixes are used to form perfective
433 partner verbs, it is traditionally assumed that the prefixes are semantically
434 empty (Avilova 1959, 1976; Čertkova 1996; Forsyth 1970; Šaxmatov 1952;
435 Švedova 1980; Tixonov 1964, 1998; Vinogradov 1972; however note that
436 some scholars have opposed this tradition, especially van Schooneveld 1958
437 and Isačenko 1960). *Gruzit'* provides an ideal testing ground for the empty
438 prefix hypothesis, since a) this verb has three supposedly empty prefixes in the
439 partner perfective verbs *zagrúzit'*, *nagrúzit'*, and *pogrúzit'* all meaning ‘load
440 (perfective)’; and b) all four verbs (imperfective *gruzit'* and all three
441 perfectives) can appear in two competing constructions, the theme – object
442 construction *gruzit' jaščiki na telegu* ‘load boxes onto the cart’, and the goal –
443 object construction *gruzit' telegu jaščikami* ‘load the cart with boxes’.

444 The point is to show that the prefixes provide different environments for the
445 constructions and because prefixes do not behave identically they are therefore
446 not identical in function or meaning. We discover that *nagruzit'* strongly
447 prefers the goal – object construction, *pogruzit'* almost exclusively prefers the
448 theme – object construction, whereas *zagruzit'* has a more balanced
449 distribution. Thus one can say that each prefix has a unique characteristic
450 preference pattern. Our analysis shows that this is a robust finding, even when
451 we take into account relevant additional environmental variation, namely the
452 use of the prefixes in constructions with passive participles, as in *Irina*
453 *Vladimirovna šla nagružennoj sumkami i sumočkami* ‘Irina Vladimirovna
454 walked along, loaded with bags and pouches’, and the use of reduced
455 constructions where one of the participants is missing, as in *mužiki gruzili les i*
456 *kamen* ‘the men loaded timber and rock’ (where the goal argument is not
457 mentioned).

458 [A description of the dataset](#)³ is provided in (1). The aim of a statistical
459 model for this dataset is to predict the CONSTRUCTION based on the
460 predictors VERB, REDUCED, and PARTICIPLE. This prediction can be
461 modeled using all three kinds of models considered here: logistic regression,
462 tree & forest, and NDL.

Gelöscht: Table 1

Gelöscht: provides a

Gelöscht: a

Gelöscht: .

464 (1) Description of the *gruzit'* dataset

- 465 • *Dataset and R script*
- 466 *datLOAD.csv; LOAD.R*
- 467 • *Source of dataset*
- 468 *Russian National Corpus (www.ruscorpora.ru)*

³ This dataset and the logistic model were presented in Sokolova, Lyashevskaya and Janda (2012).

- 473 • *Size of dataset*
 474 1920 rows (observations), each representing an example sentence containing
 475 *gruzit'*, *nagruzit'*, *zagruzit'* or *pogruzit'* 'load'
- 476 • *Rival forms*
 477 | Theme – object construction vs. goal – object construction, represented as
 478 CONSTRUCTION with values: theme, goal
- 479 • *Predictors*
 480 ○ VERB
 481 | Zero (for the unprefix verb *gruzit'*), na, za, and po
- 482 ○ REDUCED
 483 | Yes (construction is reduced) or no (full construction)
- 484 ○ PARTICIPLE
 485 | Yes (passive participle) or no (active form)
 486

Gelöscht: t

Gelöscht: z

Gelöscht: y

Gelöscht: y

487 3.1.1 Logistic regression

488 The optimal logistic regression model for this dataset includes all three
 489 predictors as main effects, plus an interaction between the verb and participle
 490 predictors. The formula for this model is (the asterisk '*' tells R to include not
 491 only VERB and PARTICIPLE as main effects, but also their interaction):⁴

493 CONSTRUCTION ~ VERB + REDUCED + PARTICIPLE + VERB*PARTICIPLE

495 | The linear model yields the estimates for the coefficients shown in Table 1.
 496 This table may seem rather daunting, but the basic ideas underlying these
 497 numbers are straightforward. The first column, labeled 'Estimate', presents the
 498 estimated coefficient. To interpret the values of the coefficients, recall that a
 499 logistic model estimates how the log of the odds ratio depends on the
 500 predictors. For an odds ratio, we need to know what R considers to be a
 501 success and what it takes to be a failure. By default, R will order the levels of

Gelöscht: 2

⁴ Note that because any predictor that is present in an interaction is also automatically considered as a main effect, this formula can be rendered more succinctly as: CONSTRUCTION ~ VERB*PARTICIPLE + REDUCED. The LOAD.R script tracks how this formula was arrived at through successive iterations, gradually increasing the number of predictors and comparing the results. Further interactions were not found to be statistically significant.

507 the response alphabetically, and take the second one to be a success. For the
 508 present data, this means that the theme construction is a success, and that the
 509 model ascertains how the log of the number of theme constructions divided by
 510 the number of goal constructions depends on the predictors.

511

512 **Table 1** Coefficients for logistic regression model of LOAD data

Gelöscht: 2

	Estimate	Std. Error	Wald Z	p-value
Intercept	-0.946	0.202	-4.679	0.0000
VERB=po	6.714	1.022	6.570	0.0000
VERB=za	1.092	0.245	4.455	0.0000
VERB=zero	2.334	0.245	9.539	0.0000
PARTICIPLE=yes	-4.186	1.022	-4.096	0.0000
REDUCED=yes	-0.889	0.175	-5.085	0.0000
VERB=po, PARTICIPLE=yes	3.895	1.598	2.438	0.0148
VERB=za, PARTICIPLE=yes	1.409	1.077	1.308	0.1910
VERB=zero, PARTICIPLE=yes	-1.772	1.441	-1.229	0.2190

513

514

515 The list of estimates for the coefficients begins at the Intercept. The way in
 516 which R by default deals with factors is to take one factor level as point of
 517 reference. For this particular factor level, e.g., no for the factor REDUCED, the
 518 group mean is calculated. For the other factor level (yes), the difference
 519 between its group mean and the group mean for no (the reference level) is
 520 calculated. All group means are on the logit scale.

521 R chooses as values at the Intercept those that come first alphabetically
 522 (unless the user specifies otherwise). Thus the Intercept here involves these
 523 values for the three predictors: VERB=na, PARTICIPLE=no, REDUCED=no.
 524 The intercept has the value -0.9465, indicating that for the subset of data for
 525 which VERB=na, PARTICIPLE=no, and REDUCED=no, the theme
 526 construction is used less often than the goal construction (the odds ratio is less
 527 than one, and the log of a number between 0 and 1 is negative). When we

529 change to another group mean, for VERB=na, PARTICIPLE=no, and
 530 REDUCED=yes, the group mean is $-0.9465 - 0.8891 = -1.8356$, indicating that
 531 for REDUCED observations, the theme construction is an even smaller
 532 minority.

533 The interpretation of VERB and PARTICIPLE requires special attention,
 534 because these two predictors enter into an interaction. The interaction
 535 introduces additional adjustments that have to be applied when the factors
 536 involved in the interaction both have values that differ from the reference
 537 values. The eight group means can be constructed from the estimates of the
 538 coefficients as follows, [see Table 2](#):

540 **Table 2** Interpreting interacting predictors

Interacting predictors	Calculation of group means
VERB=na, PARTICIPLE=no	-0.9465
VERB=po, PARTICIPLE=no	-0.9465+6.7143
VERB=za, PARTICIPLE=no	-0.9465+1.0920
VERB=zero, PARTICIPLE=no	-0.9465+2.3336
VERB=na, PARTICIPLE=yes	-0.9465-4.1862
VERB=po, PARTICIPLE=yes	-0.9465+6.7143+3.8953-4.1862
VERB=za, PARTICIPLE=yes	-0.9465+1.0920+1.4087-4.1862
VERB=zero, PARTICIPLE=yes	-0.9465+2.3336-1.7717-4.1862

541
 542 Thus, for VERB=zero, PARTICIPLE=yes, REDUCED=no, the model predicts
 543 a log odds ratio equal to -4.5708, which converts (with the plogis function) to a
 544 proportion of 0.0102. This compares well with the observed counts, 90 for goal
 545 and 1 for theme (proportion for theme: 0.0110).

546 The second column in [Table 1](#) presents a measure of how uncertain the
 547 model is about the estimate for the coefficient. The greater this measure, the
 548 standard error, the more we should be on guard. The third column is obtained
 549 by taking the values in the first column and dividing them by the values in the
 550 second column, resulting in so-called Z scores. These Z scores follow a

Gelöscht: x

Gelöscht: 2

553 standard normal distribution, and the final column with p-values presents a
554 measure of how surprised we should be that the scores are as big as they are.
555 More specifically, p-values evaluate how surprised we should be to observe a
556 coefficient with as large (or as small, when negative) a value as actually
557 observed, where we evaluate surprise against the possibility that the predictor
558 is not associated with the response at all, i.e., that the values of the predictors
559 and the response are random. The standard cutoff for recognizing statistical
560 significance in our field is $p = 0.05$, but it should be kept in mind that for large
561 datasets, and for data with much better experimental control than we usually
562 have in language studies, the cutoff-value can be set much lower. The values
563 for the first six lines in [Table 1](#) are all < 0.0001 . For the intercept, the small p-
564 value indicates that the group mean for VERB=na, REDUCED=no,
565 PARTICIPLE=no has a log odds that is significantly below 0. Translated into
566 proportions, this means that the proportion of the theme construction is
567 significantly below 50%. For the other terms with small p-values, we have
568 good evidence that the differences in group means are significant.

569 The interaction of VERB and PARTICIPLE gets lower marks, since only
570 one of the three coefficients has a p-value below 0.05. This raises the question
571 of whether the interaction is really needed. The problem here requires some
572 care. The table of coefficients ([Table 1](#)) only lists three corrections on
573 differences between group means (the interaction terms), while there are $\binom{4}{2} =$
574 6 pairwise comparisons all in all (e.g., VERB=po versus VERB=zero is
575 missing). As a consequence, we may be missing out on the most striking group
576 difference. Furthermore, when multiple coefficients are evaluated with p-
577 values, there is an increased probability of getting a low p-value by chance.

Gelöscht: the table

579 This can be corrected by applying the Bonferroni correction (Dunn 1961),
 580 which works as follows for the present example. We have 3 coefficients for the
 581 interaction, and our significance level (alpha) is 0.05. We divide alpha by 3,
 582 resulting in 0.0167. Any coefficient with a p-value less than 0.0167 is certain to
 583 be significant. So we now know that the interaction captures at least one
 584 significant contrast.

585 A second way of evaluating the interaction is to compare a model without
 586 the interaction with a model that includes the interaction. We can do this with
 587 an analysis of deviance test, which will evaluate whether the extra coefficients
 588 required for the interaction buy us a better fit to the data. In fact, we can apply
 589 this approach to a sequence of models, each one having one more predictor
 590 than the previous one. If we start with a model with just an intercept (the grand
 591 mean, model 1), and then add in first VERB, then PARTICIPLE, then
 592 REDUCED, and finally the interaction of VERB by PARTICIPLE (model 5),
 593 we obtain Table 3.

594

595 **Table 3** Model comparison statistics for the LOAD data

	Resid. Dev	Df	Deviance	p-value	Reduction in AIC
Intercept	2645.16				
Verb	1305.31	3	1339.85	0.0000	1333.8
Participle	950.73	1	354.58	0.0000	352.6
Verb:Participle	933.48	3	17.25	0.0006	11.2
Reduced	906.69	1	26.80	0.0000	24.8

596

597

598 The column named Resid. Dev lists the residual deviance, the unexplained
 599 variation in the data. As we include more predictors, the residual deviance
 600 decreases. The column labeled Df specifies how many coefficients were
 601 required to bring the residual deviance down. How much the deviance was

602 reduced is given by the column labeled Deviance. The column with p-values
603 shows that each reduction in deviance is significant. Finally, the last column
604 lists the reduction in AIC, a measure of goodness of fit that punishes models
605 for having many coefficients. The reduction in AIC accomplished by a
606 predictor is an excellent guide to its importance. Here, we see that VERB is
607 most important, followed by PARTICIPLE, followed by REDUCTION,
608 followed by the interaction of VERB by PARTICIPLE.

609 The *C* value (concordance index; this is one of the statistics yielded by the
610 logistic regression—see the R code and output on
611 ansatte.uit.no/laura.janda/RF/RF.html) of 0.96 tells us that the fit of the model
612 is excellent. The accuracy of the model is 89%, where we judge the model to
613 make a correct prediction if the estimated probability for the theme
614 construction is greater than or equal to 0.5 and the theme construction was
615 actually observed.

616

617 **3.1.2 Tree & forest**

618 The tree & forest analysis gives entirely parallel results. Here our formula is:

619

620 CONSTRUCTION ~ VERB + REDUCED + PARTICIPLE

621

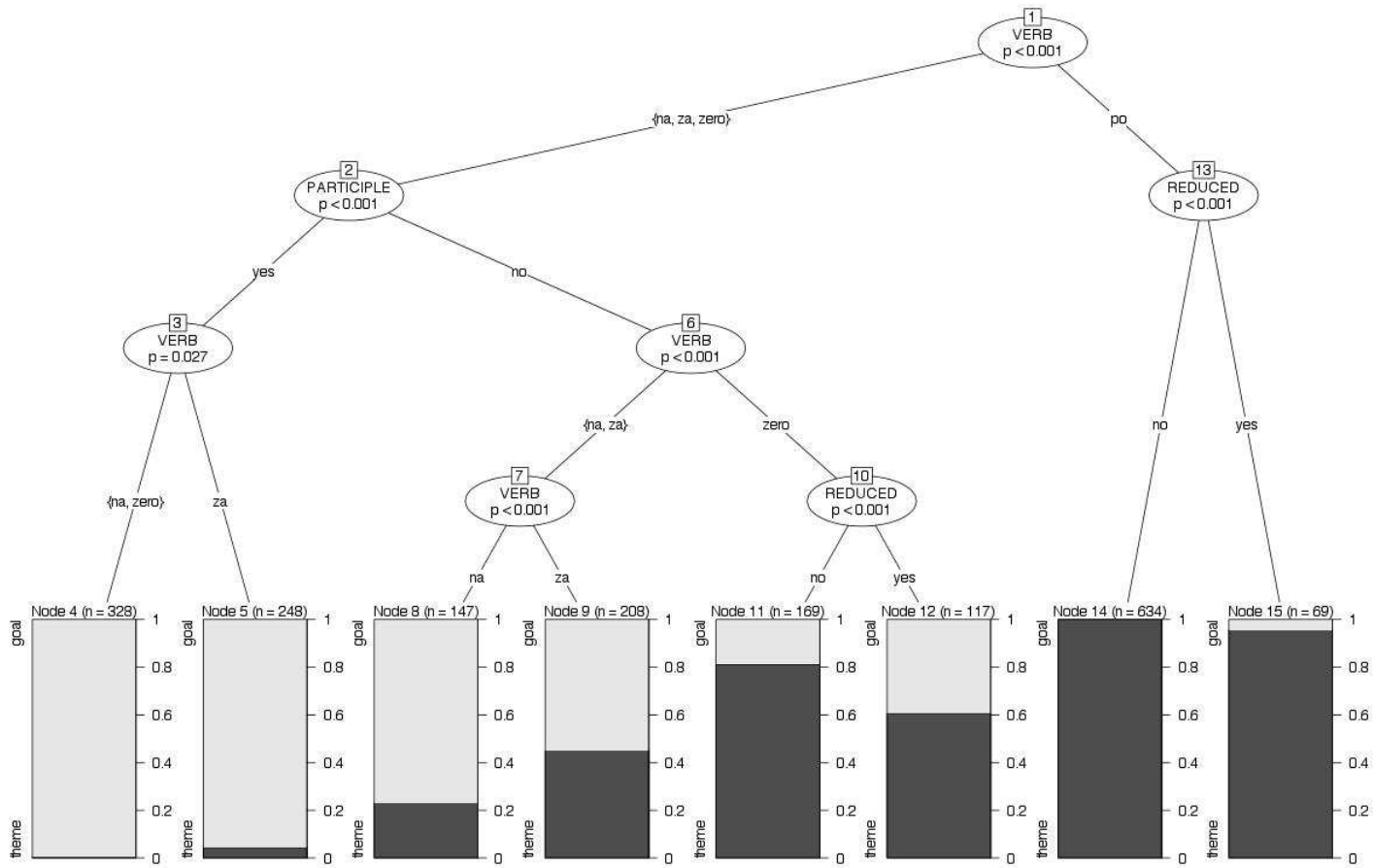
622 In the tree & forest analysis we can skip the tedium of testing different model
623 equations. We don't have to worry about how many predictors we put in, nor
624 do we have to specify interactions. Both the classification tree and the
625 classification forest will eliminate any predictors that are not significant and
626 interactions are taken into account automatically, as described below.

627 Figure 2 summarizes graphically the results of the recursive partitioning
628 tree. The first split is on VERB, distinguishing po (for which the theme is
629 almost always used) from the other three cases for which the theme is less
630 probable. The p-value in the oval presents a measure of surprise for how well
631 separable the theme and goal realizations are given information about the level
632 of VERB. The algorithm considers all possible splits, not only for VERB, but
633 also for PARTICIPLE and REDUCED, and chooses the predictor (and the
634 combination of levels of that predictor) that separates the theme and goal
635 constructions best. The choice of the best splitting criterion is made locally.
636 The algorithm does not look ahead to see whether an initial split that is not as
637 good might be offset by later greater gains. As a consequence, the predictor
638 providing the first split is often one of the most important predictors, but it is
639 not necessarily true that it is the most important predictor.

640 Once a split has been made, the same procedure (finding the locally best
641 splitting criterion, if any) is applied to both subsets (in the present case, po
642 versus na, za, zero). In this way, the dataset is recursively partitioned into
643 increasingly smaller subsets that are more homogeneous with respect to the
644 choice between theme and goal. If we go to the right branch of the tree and
645 look for the strongest factor within that branch, which is REDUCED (also with
646 $p < 0.001$), we find a split with yes on the right and no on the left. Within these
647 new subsets, further significant splits are not detected, which is not surprising
648 as choice behavior is nearly categorical here. In the left branch of the tree,
649 further splits are made on PARTICIPLE, followed by VERB and REDUCED.
650 The algorithm stops partitioning either when there is no further gain in

651 separability or when there are too few data points to allow for a meaningful
652 split.

653 The bargraph below each terminal node represents the percentage of goal
654 (light grey) vs. theme (dark grey) outcomes, and 'n=' indicates the total
655 number of datapoints in that node. So, for example, node 4 contains all of the
656 examples that involve a (past passive) participle form of either *nagruzit'* or
657 *gruzit'*; there are 328 examples of that type, and 326 (99.4%) of those have the
658 goal construction, whereas 2 (0.6%) have the theme construction. To take
659 another example, Node 9 shows us the results for active forms of *zagruzit'*:
660 there are 208 such examples, of which 114 (54.8%) have the goal construction,
661 but 94 (45.2%) have the theme construction.



662

663 Fig. 2 Recursive partitioning tree for the LOAD data

664

665 In a classification tree we see an interaction any time that the left branch of the
666 tree is different from the right branch, and/or the barplots below the terminal
667 nodes are showing different patterns. Therefore, the classification tree shows us
668 that there is in fact a complex interaction among the three factors. Within the
669 framework of a logistic regression model, one would have to include a VERB
670 by REDUCED by PARTICIPLE interaction, which would result in a large
671 number of coefficients and no noticeable improvement in goodness of fit. A
672 classification tree makes no statement about main effects, i.e., it does not
673 provide information about the effect of a given predictor with all other
674 predictors held constant. For such global statements, a logistic model should be
675 used. This having been said, it is clear that the classification tree gives us a
676 description of what is going on in the data, in a way that is visually much more
677 tractable and intuitive than the tables of figures we receive as output in the
678 regression model.

679 However, a classification tree makes its splits based on local best
680 performance, as mentioned above. Working with look-aheads would make the
681 procedure computationally intractable. In order to obtain a tree-based model
682 that avoids the risk of overfitting due to local optimization, it is useful to
683 complement the classification tree with a random forest. The random forest
684 technique constructs a large number of bootstrap samples and builds a
685 recursive partitioning tree for each of them. In order to obtain predictions from
686 this forest of trees, votes are collected from the individual trees on what they,
687 based on their training data, believe the response (e.g., goal versus theme
688 construction) to be. Typically, a random forest makes more precise predictions
689 than a standard classification tree. For the present example, the tree has a

690 classification accuracy of 88%, and the forest's accuracy increases, rather
691 atypically, only slightly to 89%. For both, $C = 0.96$.

692 The forest of trees does not provide useful information about how the
693 predictors work together. For that, we have to let ourselves be guided by the
694 classification tree. The forest does provide us with a means for assessing the
695 relative importance of the different predictors in the model. It assesses the
696 importance of a predictor, say, VERB, by randomly permuting the values of
697 VERB (na, po, za, zero) so that the relation between VERB and construction is
698 destroyed. If a predictor is truly associated with the response (theme versus
699 goal), then this procedure will cause the classification accuracy of the tree to
700 plummet. If a predictor is not predictive at all, permuting it shouldn't matter,
701 and classification accuracy should stay about the same. A measure of variable
702 importance can therefore be defined as the reduction in classification accuracy
703 under random permutation.

704 For the present data, the variable importances are 0.003 for REDUCED,
705 0.073 for PARTICIPLE, and 0.338 for VERB. VERB is the strongest predictor,
706 since a model excluding VERB is 33.8% worse than one that includes it.
707 PARTICIPLE comes next, and its removal damages the model by 7.3%. Least
708 important is REDUCED, with a value of only 0.3%. Compared to the
709 regression model, the random forest gives us comparable values for
710 concordance, with $C = 0.96$, and an accuracy of 89%.

711 Tree & forest is often an excellent choice for data with factors with few
712 factor levels. When the number of factor levels becomes large (e.g., a factor
713 VERB with 20 different verbs) and especially when there is more than one
714 factor with many factor levels, the technique becomes computationally

715 intractable. For such datasets, a mixed logistic regression model is the best
716 choice, see Sect. 3.3 for an example.

717

718 3.1.3 Naive discriminative learning

719 NDL can also be used as a classifier for the present dataset. Once again our
720 formula is simply:

721

722 CONSTRUCTION ~ VERB + REDUCED + PARTICIPLE

723

724 The NDL model yields a matrix of the weights that quantify how strongly the
725 different predictor values are associated with the rival forms goal and theme,
726 presented here in Table 4.

727 Let's see how to read this table by considering the configuration of
728 predictors VERB=no, PARTICIPLE=no and REDUCED=no. The support for
729 the theme construction is obtained simply by summing the relevant entries in
730 Table 4: $-0.25 + 0.32 + 0.22 = 0.29$. The support for the goal construction is
731 $0.45 + 0.08 + 0.18 = 0.71$. The proportional support for the theme is therefore
732 $0.29 / (0.29 + 0.71) = 0.29$. If we look at the data, we find that for this cell of
733 the design, 27 observations support the theme, and 70 the goal, i.e., 28%. This
734 fits well with the proportion predicted by NDL (29%). For any other
735 combination of predictors and their values, the calculations proceed in exactly
736 the same way.

737

738 **Table 4** NDL weights for the LOAD data

	Goal	Theme
PARTICIPALE=no	0.0794	0.3206
PARTICIPALE=yes	0.3590	0.0410

REDUCED=no	0.1757	0.2243
REDUCED=yes	0.2627	0.1373
VERB=na	0.4498	-0.2498
VERB=po	-0.4379	0.6379
VERB=za	0.3189	-0.1189
VERB=zero	0.1076	0.0924

739

740

741 From a cognitive processing perspective, the idea is that given a set of cues
742 (VERB=na, PARTICIPLE=no, REDUCED=no), activation propagates over the
743 connections of these cues with the outcomes (the goal and theme
744 constructions). The extent to which a given outcome becomes active is given
745 simply by the sum of the weights on the connections from the active cues to
746 each construction. The construction that receives most support is then the most
747 likely one to be used.

748 To assess how important a predictor is in our NDL model, we can take the
749 sum of the absolute differences of the relevant weights (for PARTICIPLE:
750 $|0.08 - 0.32| + |0.36 - 0.04| = 0.56$). The resulting values correlate extremely
751 well with the variable importance as assessed by the random forest
752 ($r = 0.9998$). Again, VERB is by far the most important factor, followed by
753 PARTICIPLE, followed by REDUCED. In other words, we get the same
754 results as in both the logistic regression and the tree & forest analyses. The
755 evaluation of the NDL model is also comparable, since it is an excellent fit
756 with $C = 0.96$ and 88% accuracy, and these figures remain unchanged under
757 ten-fold cross-validation. This example illustrates that, under ideal learning
758 conditions, human learning and statistical learning can produce nearly identical
759 results.

760 It should be noted, however, that NDL does not supply p-values of any kind.
761 It finds a set of weights that allow it to make excellent predictions given the

762 corpus data on which it is trained. For ascertaining whether a predictor is
763 statistically significant, the reader is advised to use logistic regression or a
764 classification tree.

765

766 3.2 *Pere- vs. pre-*

767 This case study addresses the question of whether the variants represent one
768 morpheme or two. *Pere-* vs. *pre-* are etymologically related prefixes, but their
769 history and behavior are quite different.⁵ In this case *pere-* is the native Russian
770 variant, whereas *pre-* is a Church Slavonic borrowing (Vasmer 1971, Vol. 3,
771 p. 356). *Pere-* has received much more attention in the scholarly literature
772 (Dobrušina, Mellina and Paillard 2001, pp. 76–80; Flier 1985; Janda 1986,
773 pp. 134–173; Shull 2003, pp. 113–119). *Pre-*, by contrast, is normally
774 mentioned only as a Church Slavonic variant (Townsend 2008, p. 59, p. 128;
775 but see Soudakoff 1975 who argues that *pere-* and *pre-* should be considered
776 distinct morphemes).

777 Our data explore variation both in terms of meaning and environment, but
778 we consistently find tendencies rather than hard-and-fast rules for the
779 distribution of forms. For example, *pere-* is usually preferred to express spatial
780 ‘transfer’, as in *perevesti* ‘lead across’, whereas *pre-* predominates in other
781 meanings such as ‘superiority’, as in *preobladat* ‘predominate’, but
782 counterexamples for this tendency are found (*preprovodit* ‘convey’ as an
783 example of a spatial ‘transfer’ use for *pre-* and *perekričat* ‘outshout’ as an
784 example of ‘superiority’ with *pere-*). In terms of environment, the most salient
785 tendencies involve a situation in which there is either prefix stacking or a +/-

⁵ Note that although these prefixes can be added to adjectives and adverbs, this case study focuses exclusively on their use with verbs.

786 shift in aspect. Prefix stacking occurs when a given verb contains more than
787 one prefix, and here *pre-* is more common, as in *prevoznesti* ‘extol’ and
788 *prepodnesti* ‘present with’, however examples with *pere-* are also found, as in
789 *pereizbrat* ‘re-elect’ and *perenaselit* ‘overpopulate’. Whereas all prefixes are
790 strongly associated with marking the perfective aspect, and thus typically serve
791 to shift the aspect of imperfective base verbs to perfective, *pre-* commonly fails
792 to effect this shift, as in *presledovat* ‘persecute’ (an imperfective verb built
793 from the imperfective base *sledovat* ‘follow’). However, *pere-* can also fail to
794 shift aspect, as in *peremenjat* ‘change’ (imperfective from imperfective base
795 verb *menjat* ‘change’),⁶ and there are also examples where both *pere-* and
796 *pre-* serve the usual role of perfectivizers, as in *pereterpet* ‘overcome’ and
797 *preterpet* ‘undergo, endure’ which are both perfective verbs from the
798 imperfective *terpet* ‘suffer’. Our analysis reveals the various strengths of the
799 semantic and environmental factors associated with *pere-* vs. *pre-* in Russian
800 verbs.

801 | [A description of the *pere-* vs. *pre-* dataset is provided in \(2\)](#). Since our goal
802 is to show that the distribution of theme – object vs. goal – object constructions
803 is affected by various factors, the aim of a statistical model for this dataset is to
804 predict the Prefix from the predictors. There are two things to note about the
805 PERE dataset that distinguish it from the LOAD dataset: 1. this data has a
806 strongly unbalanced distribution, with 1727 examples of *pere-*, but only 107
807 examples of *pre-*; and 2. this dataset includes frequency, which is a numerical,
808 quantitative predictor, as opposed to the other predictors, which are factorial

Gelöscht: Table 5 provides a

⁶ An alternative interpretation is available for this example, since *peremenjat* ‘is also the secondary imperfective of *peremenit* ‘change’.

810 (categorical, or qualitative) predictors (with discrete levels such as yes vs. no or
811 not stacked vs. stacked);

Gelöscht:

812

813 (2) [Description of the *pere-* vs. *pre-* dataset](#)

- 814 • *Dataset and R script*

815 datPERE.csv; PERE.R

- 816 • *Source of dataset*

817 Russian National Corpus (www.ruscorpora.ru)

- 818 • *Size of dataset*

819 1836 rows, each representing a verb prefixed by either *pere-* or *pre-* that is
820 attested at least once in the Russian National Corpus

- 821 • *Rival forms*

822 *pere-* vs. *pre-*, represented as Prefix with values: pere, pre

- 823 • *Predictors*

- 824 ○ *ShiftTrans*

825 [Comparison](#) of transitivity of base verb and prefixed verb, where 'intr' =
826 intransitive, 'tr' = transitive, 'no' = no existing base verb: intr-intr, intr-tr,
827 no-intr, no-tr, tr-intr, tr-tr

Gelöscht: comparison

- 828 ○ *FreqBase*

829 [Frequency](#) of the base verb in the RNC: ranges from 0 to 2694330; this
830 parameter is also available in log-transferred form as LogFreqBase.
831 Frequency distributions have long tails, and without a logarithmic
832 transformation, the highest-frequency words become atypical outliers that
833 may completely distort logistic regression models

Gelöscht: frequency

- 834 ○ *FreqPrefVerb*

835 [Frequency](#) of the prefixed verb in the RNC: ranges from 1 to 34992; this
836 parameter is also available in log-transferred form as LogFreqPrefVerb

Gelöscht: frequency

- 837 ○ *PerfectiveType*

838 [Natural](#), *spezializ*, *not applicable* (for imperfective) (cf. Janda 2007 for
839 types of perfectives)

Gelöscht: natural

- 840 ○ *SemanticGroup*

841 [Meaning](#) of the prefix (cf. Endresen forthcoming and
842 http://emptyprefixes.unit.no/pere_eng.htm): bridge, divide, interchange,
843 mix, overcome-duration, overdo, redo, seriatim, superiority, thorough,
844 transfer, transfer metaphorical, turn over, very (Note: These are the full
845 names as listed under SemanticGroupFullName; in SemanticGroup they
846 are abbreviated)

Gelöscht: meaning

847

848 3.2.1 Logistic regression

849 The optimal model for this dataset is captured by the following regression
850 equation, which has simple main effects only:

851

858 Prefix ~ ShiftTrans + PrefixStacking + ShiftAspect +
 859 SemanticGroup + LogFreqPrefVerb

860

861 This model specification yields a very large table of coefficients (see Table 5),
 862 a straightforward consequence of the large number of levels of the factor
 863 SemanticGroup. With the large number of factor levels in this dataset, the table
 864 of coefficients becomes less informative. Many of the differences in the group
 865 means for different values of ShiftAspect and SemanticGroup are not listed in
 866 the table. Two effects are easy to interpret, however. Firstly, the probability of
 867 *pre-* increases with prefixstacking, and secondly, this probability increases with
 868 the frequency of the prefixed verb: In Table 5, both predictors are paired with a
 869 positive and significant estimate.

Gelöscht: 6

Gelöscht: 6

870

871 **Table 5** Coefficients for logistic regression model of the *pre-* vs. *pre-* dataset

Gelöscht: 6

	Estimate	Std. Error	Wald Z	p-value
(Intercept)	-2.056	0.683	-3.011	0.0026
ShiftTrans=intr-tr	-0.841	0.615	-1.368	0.1712
ShiftTrans=no-intr	18.152	3540.605	0.005	0.9959
ShiftTrans=no-tr	17.103	3540.605	0.005	0.9961
ShiftTrans=tr-intr	-0.209	0.857	-0.243	0.8077
ShiftTrans=tr-tr	-0.649	0.347	-1.867	0.0619
PrefixStacking=stacked	2.755	0.490	5.620	0.0000
ShiftAspect=imp-pf	-1.485	0.409	-3.634	0.0003
ShiftAspect=no-imp	-20.160	3540.605	-0.006	0.9955
ShiftAspect=no-pf	-18.922	3540.605	-0.005	0.9957
ShiftAspect=pf-pf	-0.612	0.406	-1.507	0.1318
SemanticGroup=div	0.229	0.609	0.377	0.7062
SemanticGroup=intrch	-1.828	0.801	-2.281	0.0225
SemanticGroup=mix	-19.119	4435.633	-0.004	0.9966
SemanticGroup=ovc-dur	-0.795	0.676	-1.175	0.2402
SemanticGroup=overdo	-3.073	0.728	-4.221	0.0000
SemanticGroup=redo	-21.413	1189.419	-0.018	0.9856
SemanticGroup=seria	-19.398	1816.033	-0.011	0.9915
SemanticGroup=super	-0.110	0.690	-0.159	0.8737
SemanticGroup=thorough	-19.391	4849.044	-0.004	0.9968
SemanticGroup=transf	-2.367	0.631	-3.751	0.0002
SemanticGroup=transf-met	0.342	0.547	0.625	0.5318
SemanticGroup=turn	-19.671	5120.003	-0.004	0.9969
SemanticGroup=very	20.187	7565.807	0.003	0.9979
LogFreqPrefVerb	0.360	0.063	5.690	0.0000

872

876

877 Rather than going through all the contrasts listed in the table of coefficients, we
878 move on to assess the importance of the different predictors. In order to do this
879 we compare a sequence of nested models, beginning with a model with an
880 intercept only (the grand mean), to which we successively add the predictors
881 ShiftTrans, PrefixStacking, ShiftAspect, SemanticGroup, and
882 LogFreqPrefVerb in this order. The result is shown in Table 6, from which we
883 can read off that Semantic Group is the most important predictor, and
884 ShiftTrans the least important. The classification accuracy of this model is
885 96%, the index of concordance is $C = 0.95$.

Gelöscht: 7

886

887 **Table 6** Model comparison statistics for the *pere-* vs. *pre-* dataset

Gelöscht: 7

	Resid. Dev	Df	Deviance	p-value	AIC
Intercept	815.70				
ShiftTrans	789.17	5	26.53	0.0001	16.5
PrefixStacking	739.16	1	50.01	0.0000	48.0
ShiftAspect	694.90	4	44.26	0.0000	36.3
SemanticGroup	415.90	13	279.00	0.0000	253.0
LogFreqPrefVerb	379.56	1	36.34	0.0000	34.3

888

889

890 Interpreting the model using the table of coefficients is difficult, especially
891 because various predictors have many factor levels. One option for further
892 analysis is to simplify a predictor such as SemanticGroup, by collapsing similar
893 levels. However, often the categorization into many factor levels is well
894 motivated, and we will therefore now consider the tree & forest method, which
895 provides a simpler guide to the interpretation of the data.

896

899 3.2.2 Tree & forest

900 The formula for this analysis is nearly the same as the one for the logistic
901 regression, but it is not necessary (although not harmful either) to log-
902 transform the frequency counts for the base verb and the prefixed verb.
903 Furthermore, we have included Perfective Type as a predictor. In the logistic
904 regression, Perfective Type failed to reach significance, and we therefore do
905 not expect to see it emerge in the classification tree.

906

907
$$\text{Prefix} \sim \text{ShiftTrans} + \text{PrefixStacking} + \text{ShiftAspect} +$$

908
$$\text{PerfectiveType} + \text{SemanticGroup} + \text{FreqBase} + \text{FreqPrefVerb}$$

909

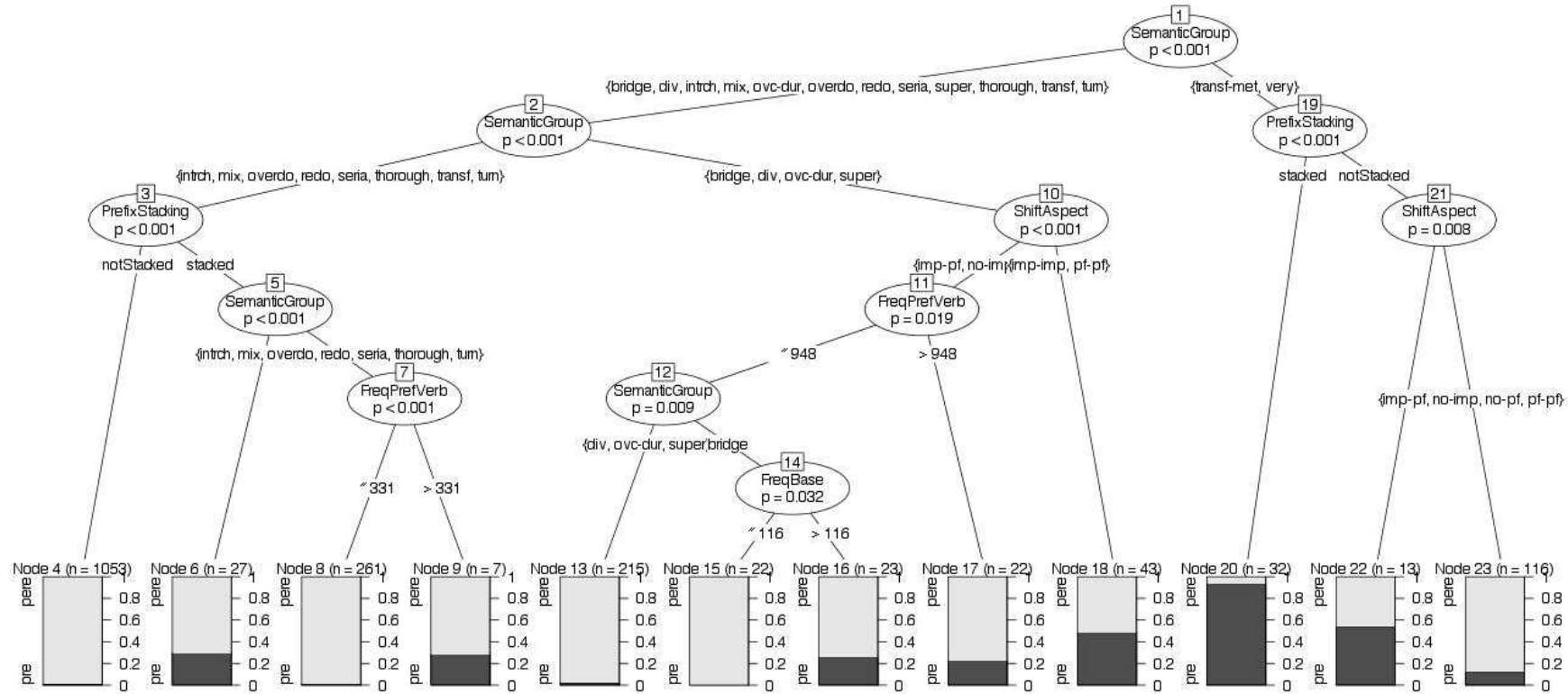
910 The recursive partitioning algorithm yields the classification tree shown in Fig.
911 3, and the random forest works out the following variable importances:
912 PerfectiveType: 0.0002, ShiftTrans: 0.0002, FreqBase: 0.0006, FreqPrefVerb:
913 0.0030, ShiftAspect: 0.0131, PrefixStacking: 0.0175, SemanticGroup: 0.0380.

914 Notice first of all that the classification tree does not include all of the
915 predictors that appear in the formula: it retains SemanticGroup, PrefixStacking,
916 ShiftAspect, FreqPrefVerb and FreqBase, but excludes ShiftTrans and
917 PerfectiveType. This fits well with the results of the logistic regression, which
918 did not support PerfectiveType at all, and which revealed ShiftTrans to be the
919 least important predictor. As promised above, the classification tree can decide
920 on its own which variables are important and which are not, and it simply
921 ignores the ones that are not important. The variable importance according to
922 the random forest is in agreement with the ranking of variable importance
923 based on the reduction in AIC for the logistic model. Interestingly, the

924 classification forest outperforms the logistic regression model: $C = 0.98$ and

925 accuracy = 96%.

926



927

928 Fig. 3 Recursive partitioning tree for the *pere-* vs. *pre-* dataset

929

930 The classification tree guides us towards a more complex interpretation of the
931 data than the logistic regression model, which only detected simple main
932 effects. From Fig. 3 it is possible to see, for instance, that for verbs from the
933 transf-met and very semantic groups, *pre-* is used almost exclusively when
934 there is no prefix stacking.

935

936 3.2.3 Naive discriminative learning

937 The observations in this dataset are a sample of the experience that the average
938 language user has with the contexts in which the choice between the rival
939 forms *pere-* vs. *pre-* arises. Therefore, NDL is an appropriate model for this
940 dataset. We are interested in whether NDL also provides a good fit to the data,
941 for two reasons. Firstly, if the model provides a good fit, it provides an
942 explanation of how language users, immersed in an environment from which
943 the corpus data are sampled, implicitly absorb and internalize the quantitative
944 forces shaping the use of *pere-* vs. *pre-*. Secondly, the the better the model fits
945 the data, the more stable we may expect the system to be.

946 The *pere-* vs. *pre-* data are especially interesting from a learning perspective
947 because these data provide information on the frequency with which forms are
948 used. In random forest and logistic regression analyses, as described above, this
949 frequency is taken into account as a property of a given data point, along with
950 other properties such as shifts in aspect or transitivity. Within the NDL
951 approach, the frequency of the derived word is not taken into account as a word
952 property, but rather as part of the learning experience. The equilibrium
953 equations that define the weights are calculated from the co-occurrence
954 frequencies of the word's properties. The frequencies of the derived words

955 codetermine these co-occurrence frequencies, and hence are taken into account
 956 for the estimation of the model's weights. Predictions of which prefix is most
 957 appropriate are derived from the weights on the links from a word's properties
 958 (such as aspect or transitivity shifting) to the prefix allomorph.

959 The model's classification performance, as estimated by the index of
 960 concordance C , is 0.97, and its accuracy is 94%. Under cross-validation, these
 961 values decrease to 0.87 and 84% respectively. It should be noted, however, that
 962 with 107 rows in the dataset (out of 1834, so 6%), which account for 16% of
 963 the occurrences of *pere-* (649757) vs. *pre-* (125668), data on *pre-* are sparse
 964 and as a consequence, crucial information about this suffix will often be lost in
 965 the training sets. Similarly, particular factor levels may not have been realized
 966 in an in-bag training set with the consequence that the model has to ignore such
 967 'unseen' factor levels altogether.

968

969 **Table 7** NDL weights for the *pere-* vs. *pre-* dataset

	Pere	Pre
PerfectiveType=natural	0.243	0.019
PerfectiveType=not-applicable	0.274	-0.012
PerfectiveType=specialized	0.025	0.238
PrefixStacking=notStacked	0.438	-0.045
PrefixStacking=stacked	0.104	0.289
SemanticGroup=bridge	0.081	-0.025
SemanticGroup=div	-0.099	0.155
SemanticGroup=intrch	0.192	-0.135
SemanticGroup=mix	0.160	-0.103
SemanticGroup=ovc-dur	0.104	-0.048
SemanticGroup=overdo	0.135	-0.079
SemanticGroup=redo	0.219	-0.163
SemanticGroup=seria	0.175	-0.119
SemanticGroup=super	-0.333	0.389
SemanticGroup=thorough	0.189	-0.133
SemanticGroup=transf	0.218	-0.162
SemanticGroup=transf-met	-0.285	0.341
SemanticGroup=turn	0.189	-0.133
SemanticGroup=very	-0.403	0.459
ShiftAspect=imp-imp	-0.153	0.310
ShiftAspect=imp-pf	0.270	-0.113
ShiftAspect=no-imp	0.013	0.144
ShiftAspect=no-pf	0.222	-0.065
ShiftAspect=pf-pf	0.190	-0.032

Kommentar [RSJ1]: Dear Laura, according to the Springer style rules, there should be a reference to Table 7 in the text. Can you make a suggestion, please?

Gelöscht: 8

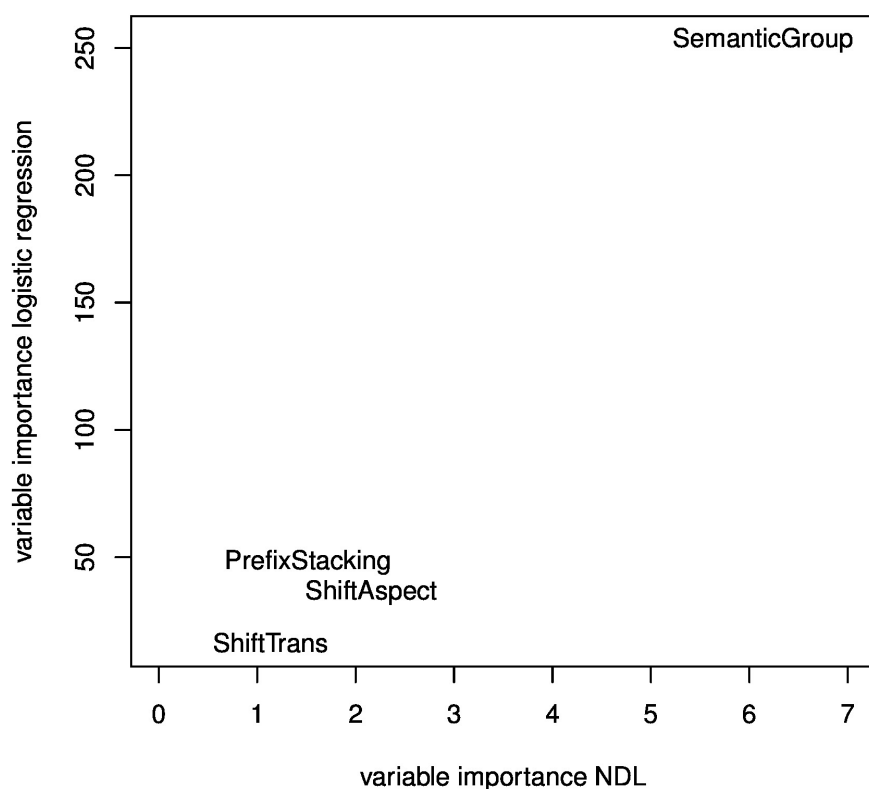
ShiftTrans=intr-intr	0.083	0.048
ShiftTrans=intr-tr	0.121	0.010
ShiftTrans=no-intr	0.135	-0.004
ShiftTrans=no-tr	0.105	0.026
ShiftTrans=tr-intr	0.002	0.129
ShiftTrans=tr-tr	0.096	0.035

971

972

973 When we assess variable importance according to NDL, we obtain the
974 following ranking: ShiftTrans: 0.55, PrefixStacking: 0.67, PerfectiveType:
975 0.72, ShiftAspect: 1.49, SemanticClass: 5.22, which hardly differs from the
976 ranking suggested by the reduction in AIC for the logistic model, as illustrated
977 in Fig. 4. What this figure shows very clearly is that the most important
978 predictor is semantic group.

979



980
 981 **Fig. 4 Variable importance according to the logistic regression model and according to**
 982 **naive discriminative learning for the *pere-* vs. *pre-* dataset**

983

984 To conclude, let's consider again how frequency of occurrence is used by the
 985 logistic regression and the classification tree on the one hand, and by NDL on
 986 the other. The logistic regression tells us that if a prefixed verb has a higher
 987 frequency, it is more likely to find *pre-* than *pere-*. This is useful information,
 988 but unless one believes that speakers have counters in their heads that keep
 989 track of how often specific forms have been used, it is information at a high
 990 level of abstraction. By contrast, the NDL model undergoes as it were the
 991 frequencies with which verbs and their distributional properties occur, and
 992 derives its predictions from the resulting discrimination weights. It is

993 conceivable, but at present far from certain, that the naive discrimination model
994 provides a cognitively more plausible assessment of the usage of *pere-* and
995 *pre-*.

996

997 **3.3 O- vs. ob-**

998 The objective of this section is to address the controversy concerning the status
999 of *o-* vs. *ob-* as either a single morpheme or two separate ones. The
1000 etymologically related variants *o-* vs. *ob-* show a complex relationship
1001 involving a variety of both semantic and phonological environments (in
1002 addition to the phonologically conditioned *obo-*). While many standard
1003 reference works (Isačenko 1960, p. 148; Timberlake 2004, p. 404; Townsend
1004 1975, p. 127; Vinogradov, Istrina and Barxudarov 1952, Vol. 1, pp. 589–592;
1005 Wade 1992, p. 277; Zaliznjak and Šmelev 1997, p. 73; Zaliznjak and Šmelev
1006 2000, p. 83), plus several specialized works (Barykina, Dobrovol'skaja and
1007 Merzon 1989; Hougaard 1973; Roberts 1981) treat *o-* and *ob-* as allomorphs of
1008 a single morpheme, some scholars (Aleksieva 1978; Andrews 1984; Krongauz
1009 1998, pp. 131–148) argue that they have split into two separate morphemes that
1010 just happen to share the same forms.

1011 The controversy is well motivated, since the behavior of *o-* vs. *ob-* covers a
1012 large portion of the space depicted in Fig. 1. We saw already in the use of
1013 *ostrič'* vs. *obstrič'* 'cut' that the two variants can sometimes be identical in
1014 terms of both meaning and environment. Additionally one can argue on the
1015 basis of examples like *okružiti* 'surround' vs. *ob'exati* 'ride around' that *o-* vs.
1016 *ob-* are classic allomorphs expressing the same meaning in phonologically
1017 complementary (non-sonorant root onset vs. sonorant root onset) environments.

1018 However, *o-* vs. *ob-* can also express a range of meanings: in addition to a
1019 meaning that can be captioned as ‘around’, as in the examples above, there are
1020 also so-called factitive uses built from adjectives meaning ‘make something be
1021 Y’ (where Y is the meaning of the base adjective or noun), as in *osložnit’*
1022 ‘make complicated’ (from *složnyj* ‘complicated’) and *obnovit’* ‘renew’ (from
1023 *novyj* ‘new’); and these two verbs additionally suggest that phonology is
1024 decisive, again with *o-* associated with a non-sonorant vs. *ob-* associated with a
1025 sonorant. However these examples give a mistaken impression: phonology is
1026 not an isolated or deciding factor, as we see in *onemečit’* ‘germanify’ (a
1027 factitive verb from *nemeckij* ‘German’) which combines *o-* with a sonorant
1028 onset, nor in *obgladit’* ‘smooth’ (a factitive verb from *gladkij* ‘smooth’) and in
1029 *obskakat’* ‘gallop around’, both of which combine *ob-* with a non-sonorant. We
1030 thus see a diverse collection of possibilities with the factors of both meaning
1031 and environment ranging from ‘same’ to various degrees of ‘different’.
1032 Additionally there is a semantic continuum between ‘around’ and the factitive
1033 type, since there are verbs like *okol’cevat’* ‘encircle’ that combine the two
1034 meanings (which can be interpreted as both a spatial sense of ‘around’ and as a
1035 factitive from *kol’co* ‘ring’). Since existing verbs and corpus data limit our
1036 opportunity to study the effects of various factors on the choice of *o-* vs. *ob-*,
1037 we present an experiment using nonce words, which give us more control over
1038 the factors. Our analysis addresses differences in meaning and differences in
1039 environment, as well as individual preferences of subjects and stems.

1040 The aim of the analysis of this dataset is to predict the choice between *o-* vs.
1041 *ob-*. There is one feature that is relevant only to part of the data: The nonce
1042 verbs were presented both as stem-stressed and as suffix-stressed, whereas the

1043 nonce adjectives were all stem-stressed. Here, we focus on the subset of the
1044 data where stress varies, i.e., the verb data.

1045 This dataset has a feature that we haven't seen in the previous analyses. In
1046 addition to comprising both quantitative (Age) and qualitative (e.g., Manner)
1047 predictors, the dataset has two predictors that have large numbers of levels:
1048 Stem (46) and Subject (60). For predictors with so many levels, it does not
1049 make sense to treat them as standard factors, which typically have only a few
1050 values which exhaustively describe all possibilities. In fact, stems and subjects
1051 are typically sampled from larger populations of stems and subjects. Under the
1052 assumption that stems and subjects are sampled randomly from these
1053 populations (an ideal that is often not truly met), these factors are referred to in
1054 the statistical literature as random-effect factors, contrasting with fixed-effect
1055 factors such as Sex (male versus female) or Voice (active, passive). Subjects
1056 and items (stems in the present example) tend to have their own specific
1057 preferences or dispreferences for a given choice (see, e.g., Dąbrowska 2008,
1058 2010; Street and Dąbrowska 2010; and Nessel, Janda and Baayen 2010, for
1059 examples from linguistics). Individual speakers, for instance, might have a
1060 personal preference for *o-* or for *ob-*. Although this dataset deals with nonce
1061 words, these nonce words will have various likenesses to real words, so we
1062 also need to weed out this potential source of extra variation in the data that
1063 could obscure the structure we are seeking to find. It will be clear that we need
1064 to bring this variability into the model in a principled way. If we fail to do so,
1065 substantial correlational structure in the model will not be accounted for, and
1066 the p-values obtained will be anti-conservative. [See below a description of the](#)
1067 [dataset in \(3\)](#):

1068

1069

(3) Description of the o- vs. ob- dataset

1070

- *Names of dataset and R script*

1071

datOB.csv; OB.R

1072

- *Source of dataset*

1073

Psycholinguistic experiment reported in Baydimirova (2010), Endresen (2011)

1074

- *Size of dataset*

1075

2630 rows, each corresponding to a response from one of sixty subjects

1076

- *Rival forms*

1077

o- vs. ob-, represented as FirstResponse1 with values: o, ob. Subjects were allowed to also make an additional response (in other words, if they first responded O, they were allowed to make a second choice of OB). We represent only the subjects' first response in this dataset.

Gelöscht:

1078

1079

1080

1081

- *Predictors*

1082

- *Subject*

1083

Anonymized subject identifier, such as A1, A2, A3, etc.

Gelöscht: a

1084

- *Stem*

1085

The nonce stem tested, such as bukl, chup, dukt, lus, etc.

Gelöscht: t

1086

- *Stimulus Type*

1087

Word class of the stimulus presented to subjects: adjective, verb

Gelöscht: w

1088

- *Onset*

1089

Onset consonant(s) of nonce stem: m, n, b, d, etc.

Gelöscht: o

1090

- *ClusterOnset*

1091

Whether the onset contained a consonant cluster: yes, no

Gelöscht: w

1092

- *PossibleWithB*

1093

Whether the Russian phonotactics allow the combination of b+ the given onset1: TRUE, FALSE. Incompatible clusters tested in the experiment are: žr, čt, žg, tk.

Gelöscht: w

1094

1095

1096

- *Place*

1097

Place of articulation of the onset: alveopalatal, dental, labial, velar

Gelöscht: p

1098

- *Manner*

1099

Manner of articulation of the onset: affricate, fricative, sonorant, stop

Gelöscht: m

1100

- *StressStimulus*

1101

Place of stress on stimulus (differentiated only for verbs; all nonce adjectives were stem-stressed): root, suffix, NotRelevant (for adjectives)

Gelöscht: p

1102

1103

- *Gender (of subject)*

1104

Male, female

Gelöscht: m

1105

- *Age (of subject)*

1106

Ranging from 18 to 59

Gelöscht: r

1107

- *EducationLevel*

1108

Higher, IncompleteHigher, Secondary

1109

- *EducationField*

1110

Humanities, Science

1111

- *SubjectGroup*

1112

Subjects were grouped according to stimulus type: A (root-stressed verb),

1113

B (suffix-stressed verb), C (root-stressed adjective)

1114

1127 Mixed-effects logistic regression makes it possible to distinguish between
1128 variability tied to subjects and items and variability linked to the predictors of
1129 primary interest. The tree & forest model, given current implementations and
1130 hardware limitations, does not scale up to data with many subjects and many
1131 items, so we will not include that model here.

1132

1133 **3.3.1 Logistic regression**

1134 In order to facilitate the interpretation of the coefficients of the model, we
1135 center Age by subtracting from each age value the mean of Age, resulting in
1136 the predictor AgeCentered. The best mixed-effects logistic model for the subset
1137 of verbs is described by the following formula:

1138

1139
$$\text{FirstResponse} \sim \text{ClusterOnset} + \text{StressStimulus} * \text{AgeCentered} +$$

1140
$$\text{Manner} + (1|\text{Stem}) + (1|\text{Subject})$$

1141

1142 The formula indicates that StressStimulus is taken into account both as a main
1143 effect and in an interaction with Age, together with a main effect of
1144 ClusterOnset. The last two terms in the formula, (1|Stem) and (1|Subject),
1145 indicate that Stem and Subject are to be treated as random-effect factors. The
1146 other predictors are treated as fixed-effect factors: they have only a fixed
1147 (usually small) number of different levels (values) that are repeatable, in the
1148 sense that one can easily build a new dataset with the same factor levels. This
1149 is not possible for subjects sampled randomly from a large population of
1150 subjects: a new random sample will contain many new subjects, and likely only
1151 subjects that have not been seen before. This explains the term ‘mixed model’:

1152 it is a model that ‘mixes’ fixed-effect and random-effect factors in one and the
 1153 same analysis (cf. Baayen 2008, Chapter 7).

1154 Table 8 lists the coefficients for the fixed-effect predictors. The intercept
 1155 represents the group mean (on the logit scale) for ClusterOnset=no,
 1156 StressStimulus=root, and Manner=affricate, for AgeCentered = 0 (which is
 1157 equivalent to Age = mean of Age), and its negative value tells us that the
 1158 model predicts *o-* here. All predictors are well-supported by low p-values,
 1159 where we should keep in mind that for Manner we see that there is one contrast
 1160 in the group means (those of sonorants and affricates) that reaches a significant
 1161 value when taking the the Bonferroni correction into account (the p-value for
 1162 this contrast is far below $0.05/3 = 0.0167$). Interestingly, when the stress is on
 1163 the suffix, the probability of using *ob-* increases with age. When the stress is on
 1164 the root, age has no such effect.

Gelöscht: 10

1165

1166 Table 8 Coefficients for a mixed-effects logistic regression model for the *o-* vs. *ob-* dataset

Gelöscht: 10

	Estimate	Std. Error	z value	p-value
(Intercept)	-0.430	0.391	-1.101	0.2710
ClusterOnset=yes	-0.596	0.236	-2.532	0.0113
StressStimulus=suffix	1.344	0.404	3.323	0.0009
AgeCentered	0.024	0.022	1.065	0.2869
Manner=fricative	0.149	0.316	0.472	0.6366
Manner=sonorant	1.079	0.348	3.104	0.0019
Manner=stop	-0.124	0.325	-0.382	0.7022
StressStimulus=suffix:AgeCentered	0.255	0.086	2.981	0.0029

1167

1168

1169 Table 9 lists the statistics for the decrease in AIC (in the column labeled AIC)
 1170 as the different terms (listed in the rows of this table) are added to the model
 1171 specification. The first row in this table compares the AIC of a model with
 1172 Subject to that of a model with only an intercept term. The large decrease in
 1173 AIC (217.6) indicates that Subject is the most important predictor. The next

Gelöscht: 11

1177 most important predictor is Stem, which comes with a reduction in AIC of
 1178 45.3. The contributions of the linguistic predictors are much smaller. It is clear
 1179 that ClusterOnset and also the interaction of StressStimulus by AgeCentered
 1180 contribute to the model fit. It is also clear that Manner is by far the most
 1181 important linguistic predictor. (The other columns in this table have the
 1182 following interpretation: logLik is the model's log likelihood, another measure
 1183 of goodness of fit. Chisq is twice the difference in logLik, which follows a chi-
 1184 squared distribution with as degrees of freedom the number of additional
 1185 parameters used by the more complex model. This number is listed in the
 1186 column labeled Chi.Df. The p-value is derived from these chi-squared
 1187 statistics.)

1188

1189 **Table 9** Model comparison statistics for the *o-* vs. *ob-* dataset

Gelöscht: 11

	logLik	Chisq	Chi.Df	p-value	Reduction in AIC
Subject	-807.13				217.6
Stem	-783.49	47	1	0.0000	45.3
ClusterOnset	-779.65	8	1	0.0056	5.7
StressStimulus	-777.96	3	1	0.0660	1.4
AgeCentered	-776.58	3	1	0.0967	0.8
StressStimulus:AgeCentered	-772.59	8	1	0.0047	6.0
Manner	-762.28	21	3	0.0001	14.6

1190

1191

1192 The index of concordance for this model is $C = 0.82$, and its accuracy is 74%.

1193

1194 3.3.2 Naive discriminative learning

1195 NDL, using the following model specification, performs equally as well as the

1196 mixed-effects model: $C = 0.82$ and an accuracy equal to 75%.

1197

1199 FirstResponse ~ ClusterOnset + StressStimulus + Age + Manner +
1200 Stem + Subject

1201

1202 It should be noted that NDL is defined only for factorial predictors. Since Age
1203 is a numerical predictor, it is automatically split on the mean into two subsets,
1204 in the present case, subjects older or younger than 24. Table [10](#) lists the
1205 weights for the main predictors, after removal of the weights for the individual
1206 stems and subjects. From this table, it is easy to see that the younger subjects
1207 prefer *o-*, whereas the older subjects prefer *ob-*. In contrast to the mixed-effects
1208 logistic regression model, the naive discrimination model supports an
1209 unconditioned effect of age. The predictors are ranked according to their
1210 measures of variable importance as follows: ClusterOnset: 0.21, Age: 0.22,
1211 StressStimulus: 0.26, Manner: 0.52, Stem: 7.66, Subject: 11.16. NDL is in
1212 agreement with the mixed-effects logistic model that Manner, Stem, and
1213 Subject are the most important predictors.

Gelöscht: 12

1214

1215 Table [10](#): NDL weights (selected) for the *o-* vs. *ob-* dataset

Gelöscht: 12

	O	Ob
Age in [18,24]	0.19	0.09
Age in [24,59]	0.07	0.20
ClusterOnset=no	0.09	0.20
ClusterOnset=yes	0.18	0.08
Manner=affricate	0.10	0.02
Manner=fricative	0.09	0.05
Manner=sonorant	-0.07	0.20
Manner=stop	0.14	0.00
StressStimulus=root	0.20	0.07
StressStimulus=suffix	0.07	0.21

1216

1217

1218 Although NDL works well for this dataset as a statistical classifier, the weights
1219 do not show a good interpretation from a learning perspective. From a

1222 cognitive perspective, it would be much more preferable to train a NDL
1223 network on the experience that speakers have with the *o-* and *ob-* rival prefixes,
1224 and then to use this network to predict what prefix speakers use for nonce
1225 verbs. In this respect, the *o-* vs. *ob-* dataset differs from the *gruzit* ‘load’ data
1226 and the *pere-* vs. *pre-* data, which comprise observations from corpora that
1227 constitute speakers’ experience with the language, and from which we can
1228 draw conclusions about what they have learned and what choices they are
1229 likely to make.

1230

1231 3.4 *nu* vs. \emptyset

1232 The objective of this case study is to chart an ongoing language change that
1233 serves to support a distinction between inchoative and stative verbs that are
1234 undergoing this change as opposed to semelfactive verbs that are not
1235 undergoing this change. Inchoative verbs such as (*ob*)*soxnut* ‘dry’ are
1236 undergoing a language change in Russian in which some past tense forms are
1237 dropping the *nu*-suffix in favor of unsuffixed (\emptyset) variants. This language
1238 change has been discussed in the scholarly literature (Bulaxovskij 1950, 1954;
1239 Černyšev 1915; Dickey 2001; Gorbačevič 1971, 1978; Nessel 1998; Plungjan
1240 2000; Rozental’ 1977; Vinogradov and Švedova 1964), but only one previous
1241 corpus study has been carried out, and that one was based on data from the
1242 1960–1970s (Graudina, Ickovič and Katlinskaja 1976, 2001, 2007). Table 11
1243 presents the relevant forms (using (*ob*)*soxnut* ‘dry’ to illustrate) and variants
1244 arranged according to overall trends identified in our case study. The left-hand
1245 side of the table presents forms for which the *nu*-variant is preferred; forms that
1246 prefer the \emptyset -variant are on the right. Vertically, each side of the table is

Gelöscht: 13

1248 ordered according to the strength of the preference, with the strongest
 1249 preference on top.

1250

1251 **Table 11** Overall preference for *nu* vs. \emptyset among inchoative and stative verbs

Gelöscht: 13

Strength of preference	Forms preferring <i>nu</i>	Forms preferring \emptyset
strongest	unprefixed participle: <i>soxnuvšij</i> > <i>soxšij</i> gerund: <i>obsoxnuv</i> > <i>obsoxšij</i>	non-masculine finite past: (<i>ob</i>) <i>soxnula</i> , - <i>o</i> , - <i>i</i> < (<i>ob</i>) <i>soxla</i> , - <i>o</i> , - <i>i</i> prefixed masculine finite past: <i>obsoxnul</i> < <i>obsox</i> unprefixed masculine finite past: <i>soxnul</i> < <i>sox</i>
weakest		

1252

1253 Since the data in this case study involves primarily inchoative and stative
 1254 verbs (plus a few transitives like *dvinut* ‘move’), there is no variation along
 1255 the meaning dimension in Fig. 1, but Table 11 gives some indication of the
 1256 complex relationships among differences in environment, since here we can
 1257 already see an interaction between the grammatical form and the presence vs.
 1258 absence of a prefix. At least two other environmental factors seem to be
 1259 involved, namely the phonological shape of the root and the presence vs.
 1260 absence of the -*sja* / -*s*’ reflexive marker. Verbs with roots ending in a velar
 1261 fricative like (*ob*)*soxnut* ‘dry’ are generally the most likely to retain *nu*,
 1262 heading a cline that proceeds through velar plosives as in (*po*)*bleknut* ‘fade’
 1263 and then dental fricatives as in (*po*)*gasnut* ‘go out’, ending with labial plosives
 1264 which are most likely to prefer \emptyset as in (*po*)*gibnut* ‘perish’. The -*sja* / -*s*’
 1265 reflexive marker also has an effect: when the marker is present, the gerund
 1266 appears in nearly equal numbers with *nu* vs. \emptyset , so forms like *proniknuvšis*’ and
 1267 *pronikšis*’, both meaning ‘having penetrated (intrans.)’ are attested
 1268 approximately equally. However, when -*sja* / -*s*’ is absent, a preference for *nu*
 1269 is maintained, so *proniknuv* is more frequent than *pronikši* ‘having penetrated

Gelöscht: 13

1272 (trans.)'. Our analysis accounts for these and additional factors along the
1273 additional diachronic dimension of change.

1274 Like the PERE dataset, NU⁷ presents us with very unbalanced data, since
1275 there are 31790 observations with \emptyset , as opposed to only 2289 with *nu*, [see \(4\)](#)
1276 [below](#). The Period and Genre predictors introduce two new types of data not
1277 present in the three datasets analyzed above, namely diachronic data and
1278 society-level data. In what follows, we focus on these two predictors.

1279

1280 [\(4\) Description of the *nu* vs. \$\emptyset\$ dataset](#)

1281 • *Name of dataset*
1282 datNU.csv1

1283 • *Source of dataset*
1284 Russian National Corpus (www.ruscorpora.ru)

1285 • *Size of dataset*
1286 34079 rows, each representing an example sentence containing an inchoative
1287 verb whose infinitive form ends in *-nut'*

1288 • *Rival forms*
1289 *nu* vs. \emptyset , represented as NU with values Nu and NoNu

1290 • *Predictors*

1291 ○ *Form (of the verb)*
1292 finite (non-masculine past tense forms), (past) gerund, mascsg (masculine
1293 past tense form), part (past active participle)

1294 ○ *Prefix*
1295 Prefixed, Unprefixed

1296 ○ *Period*
1297 1800–1849, 1850–1899, 1900–1949, 1950–1999, 2000–2010

1298 ○ *Genre*
1299 Church, fiction, massmedia, mix, nonfiction, private (as specified in the
1300 Russian National Corpus)

1301 ○ *Rootfinal*
1302 Type of root-final consonant, levels: dentalfricative, dentalplosive,
1303 labialplosive, none, velarfricative, velarplosive

1304 ○ *SemClass*
1305 Designation according stative vs. inchoative and transitive vs. intransitive,
1306 levels: InchIntr (inchoative intransitive), StatIntrans (stative intransitive),
1307 Transitive

1308 ○ *SJA*
1309 Presence vs. absence of *-sja / -s'* reflexive marker, levels: Sja, NoSja
1310

Gelöscht: c

Gelöscht: t

Gelöscht: d

Gelöscht: p

⁷ This dataset was presented in Nessel and Makarova (2011)

1315

1316 3.4.1 Logistic regression

1317 We begin with fitting a simple main effects model to the data, using the
1318 following model equation:

1319

1320
$$\text{NU} \sim \text{Form} + \text{Prefix} + \text{Genre} + \text{Rootfinal} + \text{SemClass} + \text{SJA} +$$

1321
$$\text{Period}$$

1322

1323 Table [12](#) lists the coefficients of this model. Due to the many predictors, and
1324 the many factor levels for these predictors, the number of coefficients is quite
1325 large. Most of the p-values are small, indicating that many of the listed
1326 contrasts are significant. However, the table lists only a small number of the
1327 possible comparisons of group means. For instance, for Genre, ‘church’ is the
1328 reference level, and the other genres are compared to this reference level, but
1329 not with each other.

Gelöscht: 15

1330 To quickly assess all possible comparisons of pairs, while correcting the p-
1331 values for the fact that we are performing a large number of comparisons, we
1332 can make use of the `glht` function from the `multcomp` package (Hothorn, Bretz
1333 and Westfall 2008).⁸ Figure 5 presents, for each pair of group means, the 95%
1334 confidence interval for the difference between these group means. For instance,
1335 the first row in the plot indicates that when the estimated group mean for
1336 ‘church’ is subtracted from the group mean for ‘fiction’, a 95% confidence
1337 interval (adjusted for multiple comparisons) is obtained that does not straddle
1338 zero (indicated by the vertical dashed line). From this, we can conclude that

⁸ In this example, we have made use of Tukey’s multiple comparisons method, see, e.g.,
Crawley (2002, p. 274).

1340 there is a significant difference between the two group means. Figure 5
 1341 indicates that there are two other contrasts that are significant, both involving
 1342 ‘church’. All other comparisons of pairs do not support significant differences.

1343

1344 | **Table 12.** Coefficients for the main effects logistics model for the *nu* dataset

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-5.25	0.35	-15.17	0.0000
Formgerund	8.36	0.15	55.41	0.0000
Formmascsg	2.24	0.12	18.91	0.0000
Formpart	3.98	0.12	33.22	0.0000
PrefixUnprefixed	3.08	0.11	27.21	0.0000
Genrefiction	1.04	0.32	3.23	0.0012
Genremassmedia	1.22	0.32	3.77	0.0002
Genremix	1.07	0.46	2.32	0.0203
Genrenonfiction	1.30	0.33	3.94	0.0001
Genreprivat	0.87	0.39	2.21	0.0270
Rootfinaldentalplosive	-10.17	169.96	-0.06	0.9523
Rootfinallabialplosive	-1.49	0.12	-12.58	0.0000
Rootfinalnone	-1.24	0.30	-4.10	0.0000
Rootfinalvelarfricative	-1.10	0.11	-10.22	0.0000
Rootfinalvelarplosive	-0.95	0.09	-10.36	0.0000
SemClassStatIntrans	-0.45	0.10	-4.35	0.0000
SemClassTransitive	2.07	0.09	21.81	0.0000
SJASja	-0.55	0.12	-4.54	0.0000
Period1850–1899	-0.91	0.13	-6.76	0.0000
Period1900–1949	-1.60	0.13	-12.63	0.0000
Period1950–1999	-1.97	0.13	-15.48	0.0000
Period2000–	-1.90	0.13	-14.53	0.0000

Gelöscht: 15

Gelöscht: Table of c

1345

1346 | **Table 13.** Counts of occurrences of *nu* and \emptyset , and the proportion of *Nu*, for 5 successive half-
 1347 century periods

Period	NoNu	Nu	Proportion
1800–1849	1073	239	0.182
1850–1899	3290	348	0.096
1900–1949	8012	554	0.065
1950–1999	10810	605	0.053
2000–	8605	543	0.059

Gelöscht: 16

1348

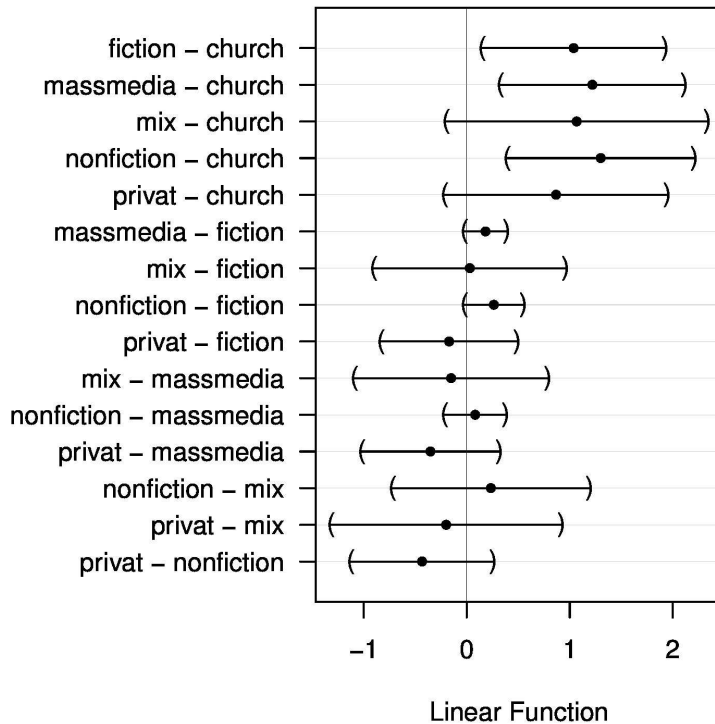
1349 Next, consider the coefficients for Period. The reference level for this factor is
 1350 1800–1849, and the four coefficients listed therefore compare later half
 1351 centuries with the first half of the nineteenth century. First note that all four
 1352 coefficients are negative. This indicates that at later moments in time, *nu* was
 1353 used less often. Also note that the coefficients become more negative as time

1357 proceeds. Only the coefficient for the most recent period, is no longer more
1358 negative than that of the preceding period. This indicates that *nu* has been used
1359 progressively less frequently over the last two hundred years, with this process
1360 of attrition possibly coming to a halt in the 21st century. Table [13](#) lists, for each
1361 half-century, the number of occurrences of *nu* and \emptyset , as well as the proportion
1362 of *nu* attestations. The proportions show exactly the same pattern as the
1363 coefficients of the logistic model, unsurprisingly. A multiple comparisons test
1364 (not shown) indicates that all comparisons of pairs for the different half-
1365 centuries are significant, with the exception of the most recent pair (1950–1999
1366 versus 2000–). The index of concordance for this model is 0.95 and its
1367 accuracy is 96.3%. A slight improvement ($C = 0.955$, accuracy = 96.6%) can
1368 be obtained by including several interactions, which increases the number of
1369 coefficients to no less than 98. As the dataset is large, the small increase in
1370 accuracy still amounts to roughly a hundred additional correct classifications.
1371 Unfortunately, the model with interactions among factors has so many
1372 interactions that it is unwieldy and thus linguistically uninterpretable.

1373

Gelöscht: 16

95% family-wise confidence level



1375 Fig. 5 Tukey's all-pair comparisons between group means for Genre
1376

1377

1378 3.4.2 Tree & forest

1379 The tree & forest method turns out to support the presence of many highly
1380 complex interactions. The classification tree shown in Fig. 6, obtained with
1381 exactly the same model specification equation as used for the logistic model,
1382 represents only the tip of the iceberg by restricting the number of splits to three
1383 levels. The tree indicates that there are two conditions in which *nu* is highly
1384 likely to be present: gerunds with no SJA and with no root final plosive, and
1385 unprefixed participles. The (full) classification tree has $C = 0.964$ and accuracy
1386 = 96.7%. This compares well with the logistic model. For an evaluation of the
1387 main trends of individual predictors, the main effects logistic model is useful,

1388 for coming to grips with the interactions, the classification tree is a good guide.
1389 It should be kept in mind, though, that for the full accuracy of the tree to be
1390 achieved, the full tree (not shown) is required. In that tree (as in the logistic
1391 model with interactions), many of the minor splits may be due to stochastic
1392 variation that comes with sampling data for inclusion in a large text corpus.

1393

1394 3.4.3 Naive discriminative learning

1395 We assess the importance of the different predictors with NDL, using the same
1396 model specification as for the logistic and tree models. This model, for which
1397 $C = 0.95$ and for which accuracy = 96.3, indicates that Form is by far the most
1398 dominant predictor, followed at a large distance by Period and Semantic Class
1399 (see Fig. 7).

1400 Accuracy can be increased by adding interactions between Form and Prefix
1401 to the model, using the following model specification:

1402

1403
$$\text{NU} \sim \text{Form} * \text{Prefix} + \text{Genre} + \text{Rootfinal} + \text{SemClass} + \text{SJA} +$$

1404
$$\text{Period}$$

1405

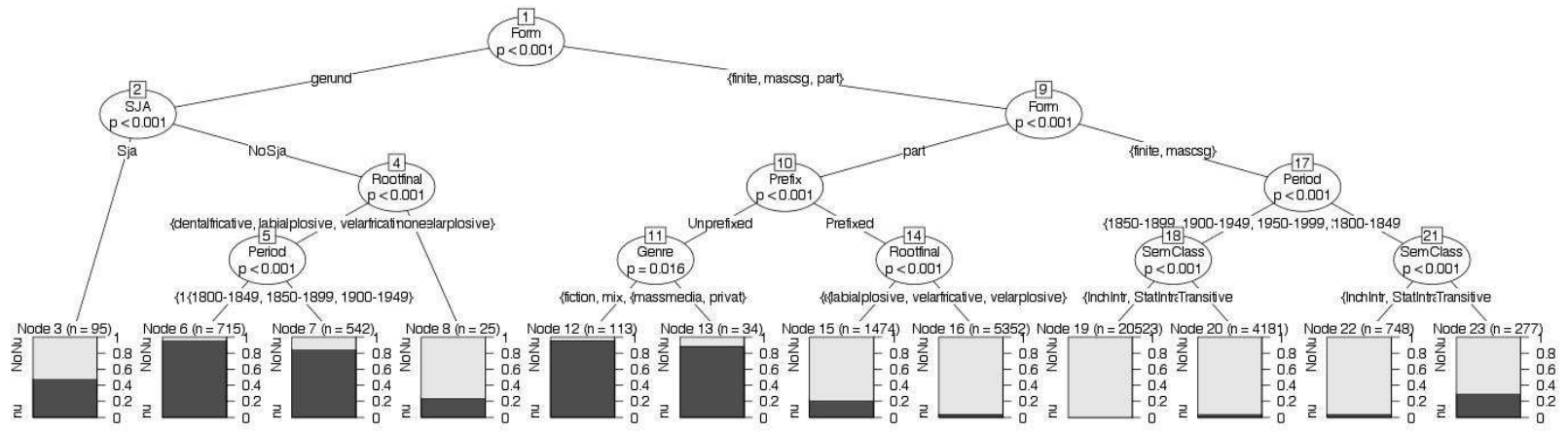
1406 This results in $C = 0.953$ and an accuracy equal to 96.7, indicating an accuracy
1407 equal to that of the other two models. The interaction asks the naive
1408 discriminative learner to add, as independent cues all unique combinations of
1409 the levels of Form and the levels of Prefix. Table 14 lists all cues and their
1410 association strengths (weights) to NoNu and Nu, ordered by the values for Nu.

1411 According to the recursive partitioning tree, the conditions favoring *nu* most
1412 were gerunds with no SJA, and unprefixated participles with no root-final

Gelöscht: 17

1414 | consonant. From Table [14](#) we can see the NDL support for these conditions,
1415 | Formgerund: $+0.326 + \text{NoSJA} + 0.089 = 0.415$ and Rootfinal none: $0.014 +$
1416 | Formpart:PrefixUnprefixed $0.432 = 0.446$. We can also clearly see that the
1417 | support for *nu* decreases over time: $0.092 \rightarrow 0.041 \rightarrow 0.016 \rightarrow 0.007 \rightarrow 0.008$.

Gelöscht: 17



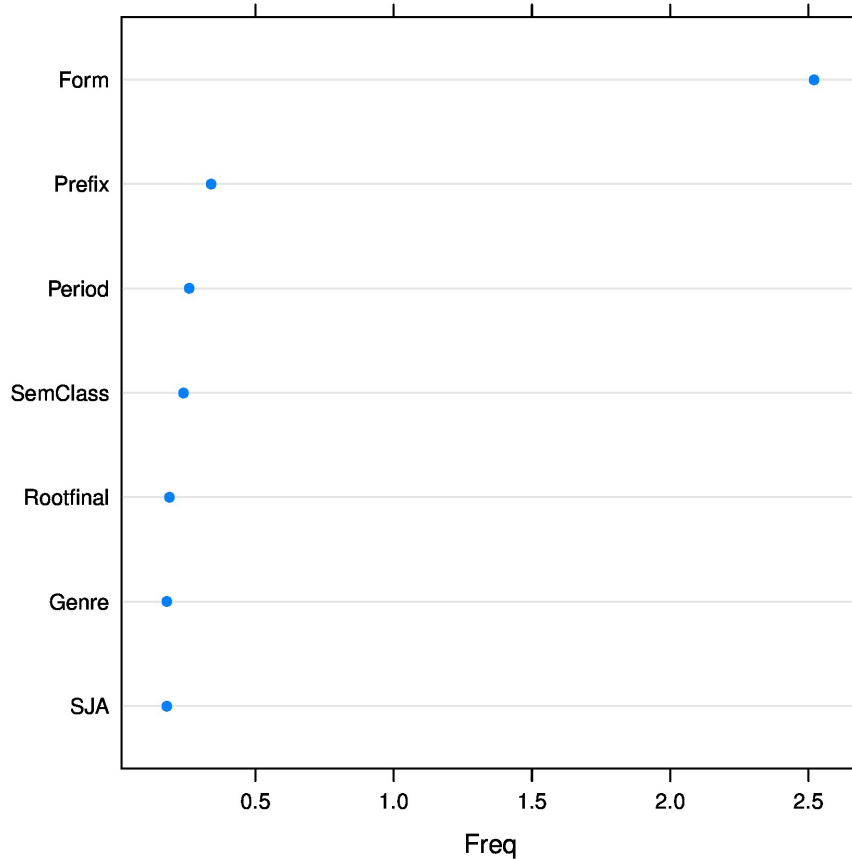
1419

1420 Fig. 6 Classification tree for the NU dataset

1421

1422

1423



1424 Fig. 7 Variable importance for the NU dataset using a simple main effects NDL model
 1425

1426

1427 4 Conclusions

1428 To conclude, we will summarize the results in two ways, firstly focusing on the
 1429 relative strengths and merits of the three statistical models used to analyze our
 1430 data and secondly interpreting the behavior of our rival forms in terms of the
 1431 relationships between their meanings and the environments they appear in.

1432

1433 Table 14 NDL weights for NoNu and Nu

Gelöscht: 17

	Weight NoNu	Weight Nu
Formpart:PrefixPrefixed	0.32	-0.28
Formfinite	0.30	-0.18
Formfinite:PrefixUnprefixed	0.24	-0.17
Formmascs	0.25	-0.13
Formmascs:PrefixUnprefixed	0.17	-0.10

Formmascsg:PrefixPrefixed	0.09	-0.04
Formfinite:PrefixPrefixed	0.07	-0.02
PrefixPrefixed	0.24	-0.01
Genrechurch	0.07	0.00
Period1950–1999	0.08	0.01
Period2000–	0.08	0.01
Rootfinallabialplosive	0.06	0.01
Rootfinalvelarfricative	0.06	0.01
Rootfinalnone	0.06	0.01
Period1900–1949	0.07	0.02
SemClassStatIntrans	0.13	0.02
Rootfinalvelarplosive	0.05	0.02
Genreprivat	0.05	0.03
Genremix	0.04	0.03
Genrefiction	0.04	0.03
Genremassmedia	0.04	0.04
SemClassInchIntr	0.11	0.04
Period1850–1899	0.05	0.04
Genrenonfiction	0.03	0.04
Rootfinaldentalfricative	0.02	0.05
Rootfinaldentalplosive	0.02	0.05
SJASja	0.15	0.07
SJANoSja	0.13	0.09
Period1800–1849	-0.00	0.09
SemClassTransitive	0.04	0.11
Formpart	-0.04	0.16
PrefixUnprefixed	0.04	0.17
Formgerund	-0.24	0.33
Formgerund:PrefixPrefixed	-0.24	0.33
Formpart:PrefixUnprefixed	-0.36	0.43

1435

1436

1437 **4.1 Pros and cons of the methods**

1438 The three statistical techniques that we have explored have different strengths
1439 and weaknesses. In what follows, we discuss these by going through a list of
1440 issues that arise in the statistical modeling of choice data.

1441

- 1442 1. Random-effect factors: The tree & forest method does not scale up for
1443 datasets with random-effect factors with many levels. We saw this for
1444 the psycholinguistic study of the distribution of *o-* vs. *ob-* in nonce
1445 words. Here, mixed-effects logistic models are the best choice.
1446 Compared to NDL, they also provide better insight into the variability
1447 associated with, for instance, speakers.

1448

1449

2. Interactions: The tree & forest method is able to detect complex interactions that are beyond the means of logistic models. The NU dataset provides an eloquent example of this. NDL can deal with complex interactions, but the weights will often not be easy to interpret.

1453

1454

3. Classification accuracy: All three techniques produce probability values for which rival form is most likely. These predictions can be used to calculate accuracy scores and indices of concordance. Across the four data sets, the different statistical methods provide very similar results, although occasionally, one method may clearly outperform the others. The general convergence, however, is reassuring, for two reasons. Firstly, it shows that we have a good understanding of the quantitative structure of the data. Secondly, we can use different methods in parallel, combining the strengths of both to compensate for individual weaknesses. For instance, a classification tree can be used to better understand interactions in a logistic model.

1465

1466

4. Variable importance: All three methods come with a method for assessing variable importance. Here too, there is remarkable convergence between methods.

1469

1470

5. P-values: Tests of significance are available for the logistic model and for the tree & forest method. Permutation tests providing p-values could

1471

1472 be added to NDL, but are currently not implemented. Therefore, NDL
1473 is not a good choice for hypothesis testing.

1474

1475 6. Cognitive interpretation: The logistic regression and the tree & forest
1476 method are statistical techniques using mathematical principles that are
1477 probably very different from those used by the brain. NDL, by contrast,
1478 is grounded in principles of human learning, and may therefore have
1479 increased cognitive plausibility, albeit still at a high level of abstraction.

1480

1481 7. Ease of interpretation: Recursive partitioning trees tend to be easy to
1482 read and provide straightforward insight into the structure of the data.
1483 However, they may become extremely complex, with many levels of
1484 branching structure, in which case interpretation becomes bewilderingly
1485 fractionated. For simple models with factors with only two or three
1486 levels, and simple interactions, the coefficients of logistic models are
1487 well-interpretable. But for more complex models, interpretation of the
1488 coefficients becomes intractable, in which case the value of the model
1489 resides in the measures of variable importance and significance tests
1490 that it provides. Interpretation will have to proceed using different
1491 means, such as cross-tabulation or recursive partitioning trees. NDL
1492 provides weights that have a simple interpretation in terms of positive
1493 (or negative) support for a rival form from a given factor level. These
1494 weights may be easier to interpret than the coefficients of a logistic
1495 model, but, as mentioned above, they do not come with p-values.

1496

1497 8. Appropriateness: All three models can be used as statistical classifiers.
1498 However, from a cognitive perspective, NDL makes sense only when
1499 the data can be viewed as a window showing a speaker's learning
1500 experience. As a consequence, it is not recommended as a model for
1501 data spanning a long time period (i.e., more than a century). Human
1502 learning is more local, and to properly model actual speakers, one
1503 would have to restrict the input data to a time interval that mirrors the
1504 average life span of a speaker.

1505

1506 9. Number of levels of response variables: Our datasets represented
1507 exclusively linguistic choices involving only two rival forms.
1508 Languages can present more complex competitions among multiple
1509 forms. However, we restricted our study in order to optimize the
1510 comparison between logistic regression (primarily designed to handle
1511 binary choices) and the tree & forest and NDL models. The latter two
1512 models can, however, be used with larger numbers of levels for the
1513 response variable. For a regression approach to datasets with a response
1514 variable with more than two levels, see Arppe (2008) and the
1515 polytomous package for R (Arppe 2012).

1516

1517 To sum up, we recommend the tree & forest method as a highly useful method
1518 complementing logistic models. Often, it will be helpful to use both in parallel.
1519 NDL is offered as an alternative that is of potential interest from a cognitive
1520 perspective. The present study is the first to show that it performs with similar
1521 accuracy as the other two methods across a variety of data samples. It is

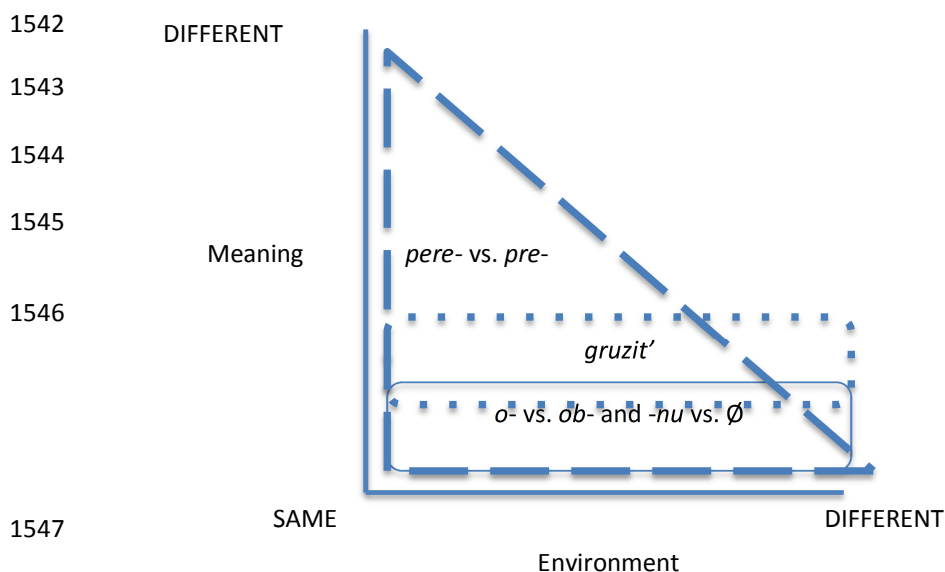
1522 conceivable that NDL may not perform as well as other methods using
1523 computational resources that are not available to the brain. By way of example,
1524 the excellent performance of random forests is due to a smart voting scheme
1525 that consults hundreds of individual trees grown on parts of the data. It seems
1526 unlikely to us that an individual's brain would work along similar lines. On the
1527 other hand, within a language group, individual speakers might be comparable
1528 to the individual trees in a forest, with the community's consensus on what
1529 form to use arising through an implicit social 'voting' scheme driven by the
1530 optimization of communication. It should therefore be kept in mind that NDL
1531 represents only low-level learning at the level of the individual, and that the
1532 forces shaping a language are much more complex. The vision behind NDL,
1533 however, is that it would be great to have a computational model that explains
1534 how grammar emerges from usage, and our current implementation should be
1535 viewed as the very first step in that direction.

1536

1537 **4.2 Rival forms and the meaning / environment plane**

1538 Where do the rival forms in our case studies fit in the space defined by
1539 variance in meaning and environment? Figure 8 gives an approximate
1540 visualization of their behavior.

1541



1548

1549 **Fig. 8 The four case studies on the meaning / environment plane**

1550

1551 For both *o-* vs. *ob-* and *nu* vs. \emptyset , only differences in environment (including

1552 both morphological and phonological environment, but also the environment of

1553 Genre for the latter) were considered while meaning was held more or less

1554 constant. The region these rival forms occupy is suggested by the thin solid line

1555 encircling '*o-* vs. *ob-* and *nu* vs. \emptyset ' in the figure. For both case studies, the rival

1556 forms can both compete in the same environment and can also be more (or

1557 less) characteristic of different environments, so they occupy a continuum

1558 between 'same' and 'different' on the bottom axis of the figure.

1559 Partially overlapping with *o-* vs. *ob-* and *nu* vs. \emptyset is *gruzit'*, represented by a

1560 dotted line. The rival forms in the *gruzit'* dataset are near-synonyms that, like

1561 the previous two sets, vary in their ability to compete in the same environments

1562 while also showing some preferences for different environments.

1563 The remaining case study is *pere-* vs. *pre-*, which is represented by a triangle

1564 with a dashed line. These rival forms cover a greater portion of the space in the

1565 figure because they can both overlap and contrast in terms of both meaning and
1566 environment.

1567 In sum, we see that different rival forms show different patterns in terms of
1568 variation in meaning and environment. This is a complicated area of linguistics
1569 that we are just beginning to explore with the help of appropriate statistical
1570 methods.

1571

1572 **References**

1573 Alekseeva, A. P. (1978). *Iz istorii pristavočnogo glagol'nogo*
1574 *slovoobrazovanija (na primere obrazovanij s ob- i o-)* (Avtoreferat
1575 kand. filol. nauk). Leningrad.

1576 Andrews, E. (1984). A semantic analysis of the Russian prepositions / preverbs
1577 O(-) and OB(-). *Slavic and East European Journal*, 28(4), 477–492.

1578 Aronoff, M. (1976). *Word formation in generative grammar* (Linguistic
1579 Inquiry Monographs, 1). Cambridge.

1580 Arppe, A. (2008). *Univariate, bivariate and multivariate methods in corpus-*
1581 *based lexicography. A study of synonymy* (PhD dissertation). University
1582 of Helsinki, Helsinki.

1583 Arppe, A. (2012). *Polytomous: Polytomous logistic regression for fixed and*
1584 *mixed effects*. R package version 0.1.4. [http://CRAN.R-](http://CRAN.R-project.org/package=polytomous)
1585 [project.org/package=polytomous](http://CRAN.R-project.org/package=polytomous).

1586 Avilova, N. S. (1959). O kategorii vida v sovremennom russkom literaturnom
1587 jazyke. *Russkij jazyk v nacional'noj škole*, 4, 21–26.

1588 Avilova, N. S. (1976). *Vid glagola i semantika glagol'nogo slova*. Moskva.

- 1589 Baayen, R. H. (2008). *Analyzing linguistic data: a practical introduction to*
1590 *statistics using R*. Cambridge.
- 1591 Baayen, R. H. (2011). Corpus linguistics and naive discriminative learning.
1592 *Brazilian Journal of Applied Linguistics*, 11, 295–328.
- 1593 Baayen, R. H., Milin, P., Đurđević, D. F., Hendrix, P., & Marelli, M. (2011).
1594 An amorphous model for morphological processing in visual
1595 comprehension based on naive discriminative learning. *Psychological*
1596 *Review*, 118(3), 438–482.
- 1597 Barykina, A. N., Dobrovol'skaja, V. V., & Merzon, S. N. (1989). *Izučeniye*
1598 *glagol'nyx pristavok*. Moskva.
- 1599 Bauer, L. (2003). *Introducing linguistic morphology*. Bristol.
- 1600 Baydimirova, A. (2010). *Russian aspectual prefixes O, OB, and OBO: A case*
1601 *study of allomorphy* (Master's thesis). University of Tromsø, Tromsø.
1602 Retrieved from <http://www.ub.uit.no/munin/handle/10037/2767>.
- 1603 Booij, G. (2005). *The grammar of words. An introduction to linguistic*
1604 *morphology*. Oxford.
- 1605 Bresnan, J., Cueni, A., Nikitina, T., & Baayen, R. H. (2007). Predicting the
1606 dative alternation. In G. Bouma, I. Kraemer, & J. Zwarts (Eds.),
1607 *Cognitive foundations of interpretation* (pp. 69–94). Amsterdam.
- 1608 Bulaxovskij, L. A. (1950). *Istoričeskij komentarij k ruskomu literaturnomu*
1609 *jazyku*. Kiev.
- 1610 Bulaxovskij, L. A. (1954). *Russkij literaturnyj jazyk pervoj poloviny XIX veka.*
1611 *Fonetika. Morofologija. Udarenie. Sintaksis*. Moskva.
- 1612 Cedergren, H. & Sankoff, D. (1974). Variable rules: Performance as a
1613 statistical reflection of competence. *Language*, 50(2), 333–355.

- 1614 Černyšev, V. I. (1915). *Pravil'nost' i čistota russkoj řeči. Opyt russkoj*
1615 *stilističeskoj grammatiki. Tom 2: Časti řeči* (2-e izd., isp. i dop.).
1616 Petrograd.
- 1617 Čertkova, M. J. (1996). *Grammatičeskaja kategorija vida v sovremennom*
1618 *russkom jazyke*. Moskva.
- 1619 Crawley, M. J. (2002). *Statistical computing. An introduction to data analysis*
1620 *using S-plus*. Chichester.
- 1621 Dąbrowska, E. (2008). The effects of frequency and neighbourhood density on
1622 adult native speakers' productivity with Polish case inflections: an
1623 empirical test of usage-based approaches to morphology. *Journal of*
1624 *Memory and Language*, 58(4), 931–951. doi:
1625 10.1016/j.jml.2007.11.005.
- 1626 Dąbrowska, E. (2010). Naive v. expert intuitions: an empirical study of
1627 acceptability judgments. *The Linguistic Review*, 27(1), 1–23. doi:
1628 10.1515/tlir.2010.001.
- 1629 Danks, D. (2003). Equilibria of the Rescorla-Wagner model. *Journal of*
1630 *Mathematical Psychology*, 47(2), 109–121. doi: 10.1016/S0022-
1631 2496(02)00016-0.
- 1632 Dickey, S. M. (2001). 'Semelfactive'- $n\phi$ and the Western Aspect Gestalt.
1633 *Journal of Slavic Linguistics*, 9(1), 25–48.
- 1634 Dobrušina, E. R., Mellina, E. A., & Pajar, D. (2001). *Russkije pristavki:*
1635 *mnogoznačnost' i semantičeskoje edinstvo. Sbornik*. Moskva.
- 1636 Dunn, O. J. (1961). Multiple comparisons among means. *Journal of the*
1637 *American Statistical Association*, 56(293), 52–64.

- 1638 Endresen, A. (Forthcoming). Allomorphy via borrowing? The status of the
1639 prefixes PRE- and PERE- in Modern Russian. Paper presented at the
1640 *7th Annual Meeting of the Slavic Linguistic Society, University of*
1641 *Kansas, Lawrence, 25–27 August 2012.*
- 1642 Endresen, A. (2011). Russkije pristavki O-, OB- i OBO-: Raznye allomorphy
1643 ili morfemy? Eksperimental'noe issledovanie. In *Proceedings of the XL*
1644 *International Philological Conference. March 14 – 19, 2011.*
1645 *Psycholinguistics* (pp. 44-55). St. Petersburg.
- 1646 Fasmer, M. (1971). *Etimologičeskij slovar' russkogo jazyka*. Vol. 3. Moskva.
- 1647 Flier, M. S. (1985). Syntagmatic constraints on the Russian prefix *pere-*. In M.
1648 S. Flier & R. D. Brecht (Eds.), *Issues in Russian morphosyntax* (UCLA
1649 Slavic Studies, 10, pp. 138–154). Columbus.
- 1650 Forsyth, J. A. (1970). *A Grammar of Aspect. Usage and meaning in the*
1651 *Russian verb*. Cambridge.
- 1652 Gorbačevič, K. S. (1971). *Izmenenie norm russkogo literaturnogo jazyka*.
1653 Leningrad.
- 1654 Gorbačevič, K. S. (1978). *VARIANTNOST' slova i jazykovaja norma. Na materiale*
1655 *sovremennogo russkogo jazyka*. Leningrad.
- 1656 Graudina, L. K., Ickovič, V. A., & Katlinskaja, L. P. (1976). *Grammatičeskaja*
1657 *pravil'nost' russkoj reči. Opyt častotno-stilističeskogo slovarja*
1658 *variantov*. Moskva.
- 1659 Graudina, L. K., Ickovič, V. A., & Katlinskaja, L. P. (2001). *Grammatičeskaja*
1660 *pravil'nost' russkoj reči. Opyt častotno-stilističeskogo slovarja*
1661 *variantov* (2-e izd., isp. i dop.). Moskva.

Kommentar [RSJ2]: oder wir geben
2012 an, weil der Vortrag in diesem Jahr
gehalten wurde.

- 1662 Graudina, L. K., Ickovič, V. A., & Katlinskaja, L. P. (2007). *Slovar'*
1663 *grammatičeskix variantov russkogo jazyka* (3-e izd., stereotipnoe).
1664 Moskva.
- 1665 Haspelmath, M. (2002). *Understanding morphology*. London.
- 1666 Hothorn, T., Bretz, F., & Westfall, P. (2008). Simultaneous inference in
1667 general parametric models. *Biometrical Journal*, 50(3), 346–363. doi:
1668 10.1002/bimj.200810425.
- 1669 Hougaard, C. (1973). Vyražat li *o/ob-* soveršaemost'? *Scando-Slavica*, 19,
1670 119–125.
- 1671 Isačenko, A. V. (1960). *Grammatičeskij stroj russkogo jazyka v sopostavlenii s*
1672 *slovackim. Morfologija*. Vol. 2. Bratislava.
- 1673 Janda, L. A. (1986). *A semantic analysis of the Russian verbal prefixes za-,*
1674 *pere-, do- and ot-* (Slavistische Beiträge, 192). München.
- 1675 Janda, L. A. (2007). Aspectual clusters of Russian verbs. *Studies in Language*,
1676 31(3), 607–648. doi: dx.doi.org/10.1075/sl.31.3.04jan.
- 1677 Krongauz, M. A. (1998). *Pristavki i glagoly v russkom jazyke: semantičeskaja*
1678 *grammatika*. Moskva.
- 1679 Matthews, P. H. (1974). *Morphology. An introduction to the theory of word-*
1680 *structure*. Cambridge.
- 1681 Nessel, T. (1998). *Russian conjugation revisited. A cognitive approach to*
1682 *aspects of Russian verb inflection* (Tromsø-Studier i Språkvitenskap /
1683 Tromsø Studies in Linguistics, 19). Oslo.
- 1684 Nessel, T., Janda, L. A., & Baayen, R. H. (2010). Capturing correlational
1685 structure in Russian paradigms: a case study in logistic mixed-effects

- 1686 modeling. *Corpus linguistics and linguistic theory*, 6(1), 29–48. doi:
1687 10.1515/CLLT.2010.002.
- 1688 Nessel, T., & Makarova, A. (2011). ‘Nu-drop’ in Russian verbs: a corpus-based
1689 investigation of morphological variation and change. *Russian*
1690 *Linguistics*, 35(4), 41–63. doi: 10.1007/s11185-011-9084-9.
- 1691 Plungjan, V. A. (2000). ‘Bystro’ v grammatike ruskogo i drugix jazykov. In L.
1692 L. Iomdin & L. P. Krysin (Eds.), *Slovo v tekste i v slovare. Sbornik*
1693 *statej k semidesjatiletiju akademika Ju. D. Apresjana* (pp. 212–223).
1694 Moskva.
- 1695 Riddle, E. M. (1985). A historical perspective on the productivity of the
1696 suffixes *-ness* and *-ity*. In J. Fisiak (Ed.), *Historical semantics.*
1697 *Historical word-formation* (Trends in Linguistic. Studies and
1698 Monographs, 29, pp. 435–461). Berlin.
- 1699 Roberts, C. B. (1981). The origins and development of O(B)- prefixed verbs in
1700 Russian with the general meaning ‘deceive’. *Russian Linguistics*, 5(3),
1701 218–233.
- 1702 Rozental’, D. È. (1977). *Praktičeskaja stilistika ruskogo jazyka*. Moskva
- 1703 Šaxmatov, A. A. (1952). *Učenie o častjax reči*. Moskva.
- 1704 Shull, S. (2003). *The experience of space. The privileged role of spatial*
1705 *prefixation in Czech and Russian* (Slavistische Beiträge, 419).
1706 München.
- 1707 Sokolova, S., Lyashevskaya, O., & Janda, L. A. (2012). The locative
1708 alternation and the Russian ‘empty’ prefixes: a case study of the verb
1709 *gruzit’* ‘load’. In D. Divjak & S. T. Gries (Eds.), *Frequency effects in*

- 1710 *language representation* (Trends in Linguistics. Studies and
1711 Monographs, 244.2, pp. 51–85). Berlin.
- 1712 Soudakoff, D. (1975). The prefixes *pere-* and *pre-*: a definition and
1713 comparison. In D. E. Davidson & R. D. Brecht (Eds.), Soviet-American
1714 Russian language contributions [Special issue]. *Slavic and East
1715 European Journal*, 19(2), 230–238.
- 1716 Street, J., & Dąbrowska, E. (2010). More individual differences in language
1717 attainment: how much do adult native speakers of English know about
1718 passives and quantifiers? In P. Hendriks & C. Koster (Eds.),
1719 Asymmetries in language acquisition [Special issue]. *Lingua*, 120(8),
1720 2080–2094. doi:10.1016/j.lingua.2010.01.004.
- 1721 Strobl, C., Malley, J., & Tutz, G. (2009). An introduction to recursive
1722 partitioning: rationale, application, and characteristics of classification
1723 and regression trees, bagging, and random forests. *Psychological
1724 Methods*, 14(4), 323–348. doi: [10.1037/a0016973](https://doi.org/10.1037/a0016973).
- 1725 Švedova, N. Ju. (Ed.) (1980). *Russkaja grammatika*. Vol. 1. Moskva.
- 1726 Tagliamonte, S. & Baayen, R. H. (2012). Models, forests, and trees of York
1727 English: *Was/were* variation as a case study for statistical practice.
1728 *Language Variation and Change*, 24(2), 135–178.
- 1729 Timberlake, A. (2004). *A reference grammar of Russian*. Cambridge.
- 1730 Tixonov, A. N. (1964). Čistovidovye pristavki v sisteme ruskogo vidovogo
1731 formoobrazovanija. *Voprosy jazykoznanija*, 1, 42–52.
- 1732 Tixonov, A. N. (1998). *Ruskij glagol. Problemy teorii i leksikografirovanija*.
1733 Moskva.
- 1734 Townsend, C. E. (1975). *Russian word-formation*. Columbus, OH.

- 1735 Townsend, C. E. (2008). *Russian word-formation*. Bloomington.
- 1736 Van Schooneveld, C. H. (1958). The so-called 'préverbe vides' and
1737 neutralization. In *Dutch contributions to the Fourth International*
1738 *Congress of Slavistics* (pp. 159–161). The Hague.
- 1739 Vinogradov, V. V. (1972). *Russkij jazyk (grammatičeskoe učenie o slove)*.
1740 Moskva.
- 1741 Vinogradov, V. V., Istrina, E. S., & Barxudarov, S. G. (Eds.) (1952).
1742 *Grammatika russkogo jazyka*. Moskva.
- 1743 Vinogradov, V. V. & Švedova, N. J. (Eds.) (1964). *Glagol, narečie, predlogi i*
1744 *sojuzy v russkom literaturnom jazyke XIX veka*. Moskva.
- 1745 Wade, T. (1992). *A comprehensive Russian grammar*. Cambridge.
- 1746 Rescorla, R. A., & Wagner, A. W. (1972). A theory of Pavlovian conditioning:
1747 variations in the effectiveness of reinforcement and nonreinforcement.
1748 In A. H. Black & W. F. Prokasy (Eds.), *Classical Conditioning II:*
1749 *current research and theory* (pp. 64–99). New York.
- 1750 Zaliznjak, A. A., & Šmelev, A. D. (1997). *Lekcii po ruskoj aspektologii*
1751 (Slavistische Beiträge, 353). München.
- 1752 Zaliznjak, A. A., & Šmelev, A. D. (2000). *Vvedenije v russkuju aspektologiju*.
1753 Moskva.
- 1754