

Inflection across the verb-noun continuum: The role of potential lexical ambiguity caused by conversion

Alexander Rauhut
Institut für Englische Philologie
Freie Universität Berlin

2021-07-15

Contents

1	Introduction	2
2	Avoiding Ambiguity	2
2.1	Conversion	3
3	Methodology	3
3.1	Data	3
3.2	Measuring conversion	3
3.3	Textual specificity	3
3.4	Lexical specificity	4
4	Analysis	4
4.1	Overview	4
4.2	Modelling inflection ratio	4
4.3	Models	5
4.4	-s vs -ed	5
5	Discussion	5
6	Outlook and Conclusion	5
	Appendix	8
6.1	Model summaries	8
6.2	Diagnostics	11
	Bibliography	13

1 Introduction

Ambiguity is pervasive in language. blablabla

Conversion is a productive process that inevitably produces homophones. Additionally, noun-verb and verb-noun conversions are doubly homophonous with inflectional affixes in *-s*.

This study will investigate potential large scale effects that the homonymy of *-s* has.

The emergence of lexical clusters is conditioned by frequency. Various studies have provided evidence (review in Diessel 2016)

In addition to the shared form, there is extensive conversion in the English language, allowing verbs to be readily used as nouns and vice-versa. However, the function of verbal *-s* and nominal *-s* is quite different, and the immediate grammatical context can be expected to be enough for disambiguation.

In section 2, I will review some of the literature on ambiguity and ambiguity avoidance and gather information from other fields of linguistics. Section 3 is an explorative overview over some ambiguous uses of *-s*. The approach taken is using POS ambiguity tags and ranked association values (log-likelihood ratio/LLR) in the same spirit as collocation analysis (Stefanowitsch & Gries 2003). In addition, I correlated the resulting association values with cosine similarities between the base forms and affixed forms taken from a simple GloVe embedding (Pennington, Socher & Manning 2014; Selivanov, Bickel & Wang 2020) trained on the respective datasets. The result shows that larger differences between base and affixed form lead to significantly fewer ambiguous tags.

Finally, in section 4, I will present a series of GAMLSS models (Generalized Additive Models for Location, Scale and Shape (Stasinopoulos, Rigby & Bastiani 2018; Rigby & Stasinopoulos 2005)) that attempt to model the influence of verb-noun or noun-verb conversion on the distribution of *-s* and *-ed* suffixation while controlling for textual dispersion, regularity of occurrence and contextual flexibility.¹

2 Avoiding Ambiguity

formal equivalence homonymy and polysemy

that deletion is not permitted in non-extraposed subject clauses Wasow (2015)

This paper demonstrates that the differences in duration of English final S as a function of the morphological function it expresses.

Hypothesis: third person singular *-s* is avoided in noun-verb conversion

In fact, it is difficult to come up with genuine examples of structural ambiguity caused by the homonymy of *-s*.

(1) The eye of the artist is concentrated on his pencil ; the pencil and the line dreams. (BNC: 49864)

Studies like Plag, Homann & Kunter (2017) call into question whether there is perfect homonymy between different types of word final *-s*. The study found significant length differences, and hypothesizes that Yung Song et al. (2013) did not find any such effect, however, albeit on a more lexically restricted study.

In a more recent study, Tomaschek et al. (2021) show that length of final *-s* can be modelled as having a discriminatory function depending on the lexical and phonological context. In other words, the duration was found to decrease with increasing contextual ambiguity. “Energy is not invested in a signal that creates confusion instead of clarity.” (Tomaschek et al. 2021: 154) This points to some degree of ambiguity sensitivity, even though it is not entirely clear what it means for ambiguity avoidance.

Systematic phonetic differences potentially contribute to the disambiguation

¹The data was preprocessed using the *Corpus Workbench* (Evert & Hardie 2011), all analyses were coded in *R* 4.1.0 (R Core Team 2021; Dowle & Srinivasan 2021), figures were created with *ggplot2* (Wickham 2016).

Wasow (2015) surveys a variety of studies ... and concludes that ambiguity avoidance has not been shown to be as common as expected considering the pervasiveness of ambiguities in language. Grice underestimated the pervasiveness of ambiguity. Provides an attempt at a taxonomy. - word order freezing as ambiguity avoidance - that deletion only in unambiguous syntactic context - German free syntax but default SVO reading. Wasow does not call it default.

Piantadosi, Tily & Gibson (2012) Some degree of ambiguity is required in an efficient communicative system. ease of processing blabla

Trott & Bergen (2020): Linguistic forms are expected to show a certain degree of ambiguity in order to provide an efficient system. Simulations suggest that ambiguity caused by homophony is more common in natural languages than expected, even when taking into account. Homophones are smoothed out in lexically neighboring areas.

2.1 Conversion

Conversion is fairly common in the English language.

3 Methodology

3.1 Data

Two datasets were used. As a general corpus, the traditional BNC (The British National Corpus 2007), and as a more specific text type, the spoken BNC2014 (Love et al. 2017) was chosen for comparison. The basic unit of analysis was lemmatized tokens, in combination with the CLAWS pos tags. It should be noted that lemmatization and pos tagging are not equivalent across the two datasets, but the results should be similar enough to have little overall effect on the following analyses. Larger samples might be desirable.

The dataset was finally broken up into verbs and nouns based on what they occur as more often. In the rare case of a tie, the item was included in both data sets. All metrics, however, were calculated on all occurrences of a lemma rather than only occurrences tagged as noun or verb as it is typically done. Like that, it is possible to capture a wider range of statistics per item, while also minimizing any a priori categorization of items, acknowledging the fact that lexical items in English are rather flexible when it comes to open word classes.

Textual co-occurrence can only serve as an upper bound. Prosodic markers in spoken language can further disambiguate most types of ambiguity. The same is true for more general types of background knowledge. In general, it can be expected that there tends to be enough structural and contextual information to disambiguate -s.

3.2 Measuring conversion

from -1 for verbal uses to +1 for nominal uses and 0 for 50/50 show LLR graph to show that it is continuous and monotonic

Proper names were commonly mistagged as verbs in the datasets, causing heavy tails in the residuals, therefore excluded from the final analysis.

3.3 Textual specificity

In order to measure the textual specificity of lexemes, and penalize forms with very specific contextually bound uses, two types of dispersion are used. The first is the Deviation of Proportions across corpus parts (DP.norm) (Gries 2008), more specifically the normalized version (cf. Lijffijt & Gries 2012). As the basic unit, the individual texts were used by text ID. DP.norm is a corpus-part-based measure bound between 0 and 1. Values close to 1 indicate a high deviation, therefore a low dispersion. It can be interpreted as a measure of how evenly tokens are spread over the corpus parts.

DP cannot account for short bursts of occurrences, therefore, Word Growth Dispersion (DWG) (Zimmermann 2020) is used in addition, which is a whole-corpus measure based on distances between occurrences. DWG is a measure of how regularly a token occurs across an entire corpus.

If ambiguity is regularly avoided, it should be expected that with growing word class ambiguity, the use of *-s* becomes more context-dependent. Considering a given dispersion or “clumpiness” of contexts in which a lexeme is likely to occur, the existence of avoidance-contexts should cause a penalty to dispersion and make the distribution clumpier.

Both measures were designed to measure dispersion or commonness of lexical items as extension to the most commonly used plain frequencies. The two measures highlight different aspects both conceptually and empirically (cf. Gries 2021 on using multiple measures of dispersion and association).

Problem: The contexts in which ambiguity is avoided/not avoided might be rather evenly dispersed themselves. This could mask potential clumpiness of *-s* occurrences.

It is difficult to identify avoidance contexts formally. Lexical and syntactical correlates of the avoided structure might be avoided as a side effect, and cause fuzzier overall differences in structure.

3.4 Lexical specificity

The ratio of hapaxes (henceforth α_1 Evert 2005: 130) on either side of a given lexeme was used as a simple measure of how fixed the immediate lexico-grammatical context is. The term *hapax* is used a bit more liberally here and refers to types that occur only once in the given window. For ease of interpretation and in order to capture the fixedness in smaller units of text, the window was held at 1 token to the left and right. Preliminary experiments with larger windows were inconclusive, so the simplest version was used. The left and right contexts were also kept separate since lumping them together conflates quite different pieces of information, considering the direction of processing, the branching structure of English, etc. For example, a token that is always preceded by a definite article as a part of a name, scores an extremely low α_1 value.

Lemmas that had *-s* forms that only occurred in totally fixed contexts, therefore scoring alpha (Hapax) values of lower than 0.3 were excluded as outliers. Exclusion of those observations improved distributional properties of the models. Examples for those outliers are abbreviations (e.g. *e.g.*, *pp.*, *etc.*), and parts of multi-word names.

On the lower frequency bands this is due to a substantial decrease of noise as it is difficult to estimate the proportion of *-s* only based on a few occurrences of *-s*. For the higher frequency bands, it is in many cases equal to the exclusion of mostly names and idiomatic expressions. This is conceptually unproblematic since multi-word structures should be treated as individual units, and as a different category of lexeme.

4 Analysis

4.1 Overview

4.2 Modelling inflection ratio

While controlling for clumpiness with dwg did not substantially change the effect of conversion, it did improve the fit of the model.

why gam: (Wood et al. 2016; Wood 2017) Semi-parametric Both extremely high and extremely low values of dwg suggest special values. Most of the variables used in the model are not expected to have a monotonic relationship with the proportion of *-s* occurrences. Very badly dispersed lexemes are overrepresented in terms of frequency, and an extremely even dispersion suggests uses typical of function words. The same is true for productivity or contextual fixedness. Both extremely high and extremely low values are untypical.

Predicting counts of *-s* directly with poisson or negative binomial regression leads to problematic model properties such as highly skewed residuals. Human perception has been found to be more sensitive to proportional changes of stimuli rather than absolute ones (cf. Kromer 2003). With that in mind, the

independent variable to be predicted was set to the proportion of occurrences of *-s*. As such the possible values can range from 0 to 1 inclusively. Therefore, the distribution chosen to be fitted was a zero-one-inflated beta distribution (Ospina & Ferrari 2012; Rigby et al. 2019). To allow for the necessary flexibility in the distribution, the model framework chosen was GAMLSS that allow the use of smoothing spline and fitting the required inflated beta distribution. The models were created with the R package *gamlss* (Rigby & Stasinopoulos 2005).

4.3 Models

4.3.1 Model fit

Table 1 shows the statistics for the residual distribution

Table 1: Statistics for residuals and fit of the three GAMLSS models

Metric	verbal -s	nominal -s	ed
mean	-0.0401	-0.0360	0.0266
variance	0.9903	0.9237	0.9075
coef. of skewness	0.4938	0.4793	-0.2081
coef. of kurtosis	5.0987	3.1296	3.6438
Filliben correlation coefficient	0.9912	0.9916	0.9956
pseudo R ²	0.6648	0.3955	0.3859

Means and variances are very close to the desired values. The model for verbal *-s* also shows a heavy right tail, which is caused by a high amount of extreme values. Arguably, this is not unusual for language data, and might be related to the sample size. All models show some degree of skew. An inherent property of the dataset is a potential for multi-modality since the basic unit of analysis is lemmas. Clustering techniques or fitting model mixes might be strategies for improvement, without making arbitrary assumptions on lexical classes and creating cut-off points, e.g. for dominantly nominal vs. dominantly verbal lexemes.

The overall deviance explained by the models is rather high for the verbal *-s* model and medium for the other two models. There is much more variation on the nominal side of the spectrum, hence more variation when it comes to plural *-s*, which is also the most frequent inflection in comparison.

4.4 *-s* vs *-ed*

Figure 1 shows a side by side comparison of the estimates for conversion. The shape is similar, but the variability in the case of nominal *s* is much larger. In both cases, there is a depression towards the middle of the continuum rather than a simple monotonic decrease. In comparison, the influence of conversion is more linear for *-ed* as can be seen in figure 2.

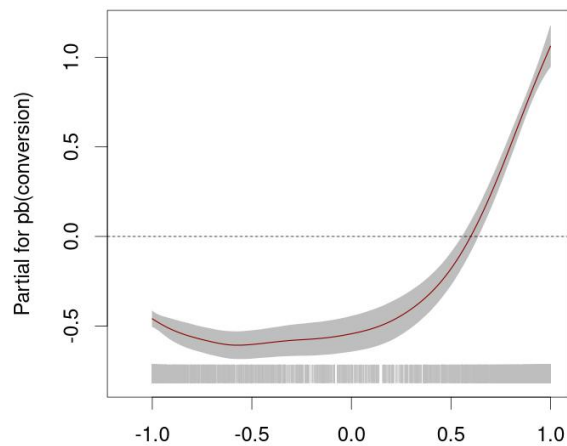
Plain, relative or log frequencies did not improve the models. In fact, model diagnostics became much worse in some cases. Possible reasons are the problematic distributional properties of word frequencies and low frequency noise. Furthermore, frequency influences every of the remaining metrics, so it is in a sense encoded there. Therefore, it was dropped in the final models presented here.

5 Discussion

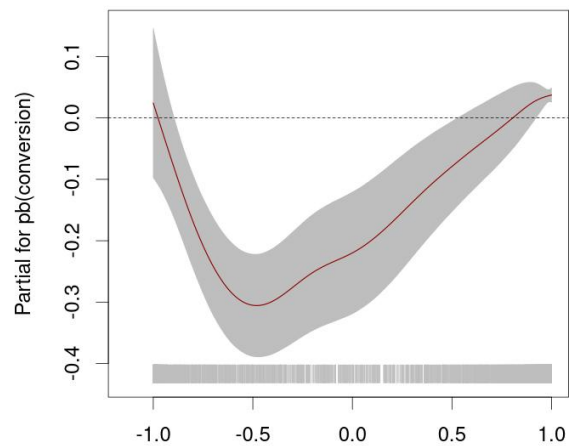
bla

6 Outlook and Conclusion

Modelling the proportion of *-s* occurrences showed a negative effect for items that are more likely to occur in conversion. This pattern was more pronounced for the nominal side of the continuum. Additionally, the



(a) verbal *-s*



(b) nominal *-s*

Figure 1: GAMLSS estimates for conversion

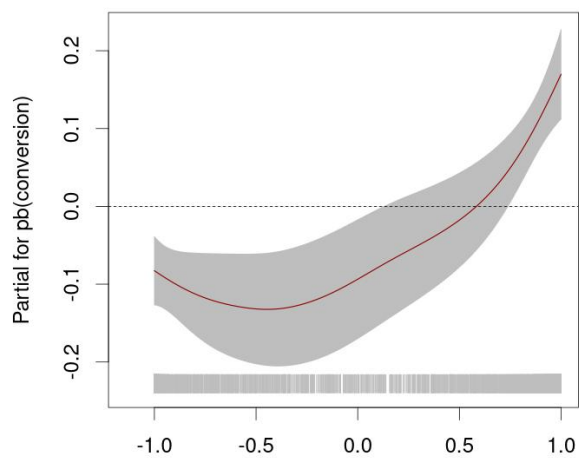


Figure 2: GAMLSS estimate for conversion for *-ed*

probability that items never occur with the *-s* suffix of their respective dominant word class was shown to drastically decrease already at rather low proportions of conversion. However, the predictions of the model have to be taken with a grain of salt, since the distributional properties of the data could not be fitted perfectly, resulting in skewed residuals. Some potential factors could be the non-randomness of the data, strong systematic noise, such as homography, and potential unaccounted multimodality, e.g. caused by other word classes. Nevertheless, the applied dispersion and specificity measures allowed to improve the fit drastically, and show promise for future improvements. Additionally, understanding morphological data as proportional, rather than count data, allows for intuitive and conceptually interesting interpretations.

-ed showed a monotonically increasing effect towards the nominal side of the spectrum indicating that nouns, if they are used as verbs, are slightly more likely to be marked with the *-ed* suffix. This effect is too slight to be conclusive, but it in comparison to *-s* there is no significant dip in the middle of the continuum. This should be expected if ambiguity plays a role.

Word embeddings and transformer networks like BERT may provide an interesting route for further studies in ambiguity (e.g. Beekhuizen, Armstrong & Stevenson 2021; Du, Qi & Sun 2019; Wiedemann et al. 2019). Clustering techniques can be used to detect homonymy (Lee 2021), which could be used in corpus annotation as a replacement of classical lemmatization. If robust and carefully applied, this could potentially lead to decreases in noise. Those recent successes in practical application are promising for the use in more descriptive application. They can provide another angle on co-occurrence statistics, and were only sparsely used in this study since there has been little systematic application in corpus linguistics. Include more contextual and “world knowledge” information, such as images (e.g. Kottur et al. 2016; Shahmohammadi, Lensch & Baayen 2021). Recent advances in context embedding might provide tools to further test where real ambiguity exists and where it affects the system of language.

More sophisticated measures of lexical fixedness/productivity than the proposed fixed-window hapax ratio are also desirable, especially to capture constructions, constructional idioms and other semi-fixed structures.

There are no significant differences in spread and frequency of *-s* in lexical forms that occur as verb-noun or noun-verb conversions. (???)

Appendix

6.1 Model summaries

6.1.1 Verbal -s

Fixed Effects

Parameter	Coefficient	SE	95% CI	t(5799.67)	p
(Intercept)	-4.79	0.24	[-5.26, -4.33]	-20.21	< .001
conversion	0.56	0.02	[0.51, 0.60]	24.16	< .001
dwg	0.08	0.37	[-0.65, 0.82]	0.22	0.827
cos_sim_s	0.90	0.05	[0.80, 1.01]	16.75	< .001
alpha1_right	1.90	0.23	[1.45, 2.35]	8.30	< .001
alpha1_left	1.47	0.20	[1.09, 1.86]	7.51	< .001

Sigma

Parameter	Coefficient	SE	95% CI	t(5799.67)	p
(Intercept)	-2.94	0.24	[-3.40, -2.47]	-12.35	< .001
conversion	0.64	0.02	[0.60, 0.68]	30.04	< .001
dwg	2.06	0.36	[1.35, 2.77]	5.67	< .001
cos_sim_s	0.30	0.05	[0.20, 0.41]	5.57	< .001
alpha1_right	0.89	0.24	[0.42, 1.36]	3.74	< .001
alpha1_left	1.17	0.21	[0.76, 1.58]	5.62	< .001

Nu

Parameter	Coefficient	SE	95% CI	t(5799.67)	p
(Intercept)	-3.37	0.75	[-4.84, -1.90]	-4.49	< .001
conversion	3.12	0.09	[2.93, 3.30]	33.10	< .001
dwg	0.34	1.00	[-1.62, 2.30]	0.34	0.733
cos_sim_s	-2.29	0.18	[-2.64, -1.95]	-12.94	< .001
alpha1_right	1.31	0.73	[-0.13, 2.74]	1.78	0.075
alpha1_left	1.79	0.71	[0.40, 3.17]	2.52	0.012

Tau

Parameter	Coefficient	SE	95% CI	t(5799.67)	p
(Intercept)	-14.90	3.07	[-20.92, -8.88]	-4.85	< .001
conversion	5.21	2.01	[1.27, 9.16]	2.59	0.010
dwg	-1.28	2.52	[-6.21, 3.66]	-0.51	0.613
cos_sim_s	-2.31	0.64	[-3.58, -1.05]	-3.59	< .001
alpha1_right	2.28	1.88	[-1.40, 5.97]	1.21	0.225
alpha1_left	9.68	2.68	[4.42, 14.94]	3.61	< .001

No. of observations in the fit: 5923
 Degrees of Freedom for the fit: 123.3256
 Residual Deg. of Freedom: 5799.674
 at cycle: 18

Global Deviance: -9870.102
 AIC: -9623.451
 SBC: -8798.822

6.1.2 Nominal -s

Fixed Effects

Parameter	Coefficient	SE	95% CI	t(12082.34)	p
(Intercept)	-4.61	0.15	[-4.91, -4.31]	-29.90	< .001
conversion	0.12	0.02	[0.08, 0.16]	6.10	< .001
dwg	1.35	0.20	[0.96, 1.74]	6.80	< .001
cos_sim_s	1.37	0.04	[1.30, 1.44]	36.58	< .001
alpha1_right	1.59	0.15	[1.28, 1.89]	10.24	< .001
alpha1_left	2.68	0.15	[2.39, 2.97]	18.36	< .001

Sigma

Parameter	Coefficient	SE	95% CI	t(12082.34)	p
(Intercept)	0.45	0.14	[0.17, 0.72]	3.18	0.001
conversion	0.03	0.02	[-0.01, 0.06]	1.39	0.165
dwg	0.25	0.18	[-0.11, 0.60]	1.35	0.176
cos_sim_s	-0.91	0.04	[-0.98, -0.84]	-25.25	< .001
alpha1_right	-0.74	0.14	[-1.01, -0.46]	-5.16	< .001
alpha1_left	0.15	0.13	[-0.11, 0.42]	1.15	0.252

Nu

Parameter	Coefficient	SE	95% CI	t(12082.34)	p
(Intercept)	-1.32	1.30	[-3.86, 1.22]	-1.02	0.308
conversion	-3.10	0.16	[-3.41, -2.79]	-19.53	< .001
dwg	-4.92	2.09	[-9.02, -0.82]	-2.35	0.019
cos_sim_s	-1.98	0.32	[-2.60, -1.36]	-6.25	< .001
alpha1_right	-2.54	1.24	[-4.96, -0.11]	-2.05	0.040
alpha1_left	1.20	1.00	[-0.75, 3.15]	1.20	0.229

Tau

Parameter	Coefficient	SE	95% CI	t(12082.34)	p
(Intercept)	-8.10	1.41	[-10.86, -5.34]	-5.76	< .001
conversion	-1.32	0.12	[-1.55, -1.09]	-11.13	< .001
dwg	4.55	1.75	[1.13, 7.97]	2.61	0.009
cos_sim_s	-1.16	0.40	[-1.95, -0.36]	-2.86	0.004
alpha1_right	0.01	1.36	[-2.65, 2.67]	9.10e-03	0.993
alpha1_left	4.70	1.31	[2.14, 7.26]	3.59	< .001

```
-----
No. of observations in the fit: 12198
Degrees of Freedom for the fit: 115.6586
Residual Deg. of Freedom: 12082.34
                           at cycle: 16
```

```
Global Deviance: -10125.98
AIC: -9894.66
SBC: -9037.742
-----
```

6.1.3 -ed

Fixed Effects

Parameter	Coefficient	SE	95% CI	t(5859.74)	p
(Intercept)	-1.60	0.24	[-2.08, -1.13]	-6.59	< .001
conversion	0.11	0.02	[0.07, 0.14]	5.41	< .001
dwg	0.17	0.39	[-0.59, 0.92]	0.44	0.661
alpha1_right	1.85	0.23	[1.39, 2.30]	7.92	< .001
alpha1_left	0.16	0.20	[-0.24, 0.56]	0.80	0.423

Sigma

Parameter	Coefficient	SE	95% CI	t(5859.74)	p
(Intercept)	1.08	0.23	[0.64, 1.53]	4.75	< .001
conversion	0.02	0.02	[-0.02, 0.06]	1.00	0.317
dwg	0.16	0.37	[-0.55, 0.88]	0.44	0.658
alpha1_right	-1.39	0.23	[-1.84, -0.94]	-6.08	< .001
alpha1_left	-0.56	0.20	[-0.96, -0.16]	-2.72	0.007

Nu

Parameter	Coefficient	SE	95% CI	t(5859.74)	p
(Intercept)	-9.41	0.64	[-10.67, -8.15]	-14.60	< .001
conversion	1.70	0.07	[1.57, 1.83]	25.83	< .001
dwg	2.91	0.90	[1.14, 4.68]	3.22	0.001
alpha1_right	2.96	0.64	[1.72, 4.21]	4.66	< .001
alpha1_left	7.21	0.62	[6.00, 8.43]	11.65	< .001

Tau

Parameter	Coefficient	SE	95% CI	t(5859.74)	p
(Intercept)	-11.33	1.10	[-13.48, -9.18]	-10.33	< .001
conversion	3.39	0.37	[2.66, 4.11]	9.17	< .001
dwg	0.41	1.50	[-2.53, 3.35]	0.27	0.783
alpha1_right	3.14	1.12	[0.95, 5.33]	2.81	0.005
alpha1_left	6.62	1.09	[4.47, 8.76]	6.05	< .001

No. of observations in the fit: 5923
Degrees of Freedom for the fit: 63.2646
Residual Deg. of Freedom: 5859.735
at cycle: 7

Global Deviance: 4295.006
AIC: 4421.535
SBC: 4844.56

6.2 Diagnostics

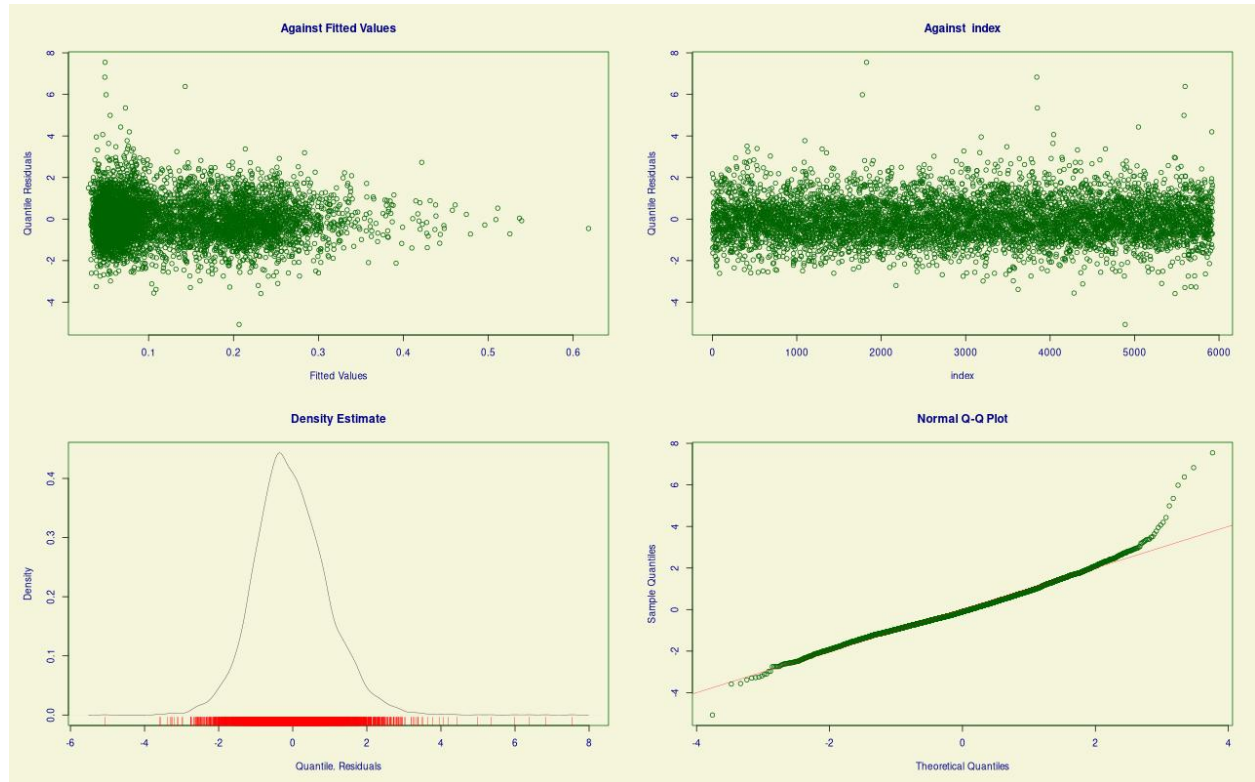


Figure 3: Diagnostics for beta inflated GAMLSS model predicting proportions of verbal -s

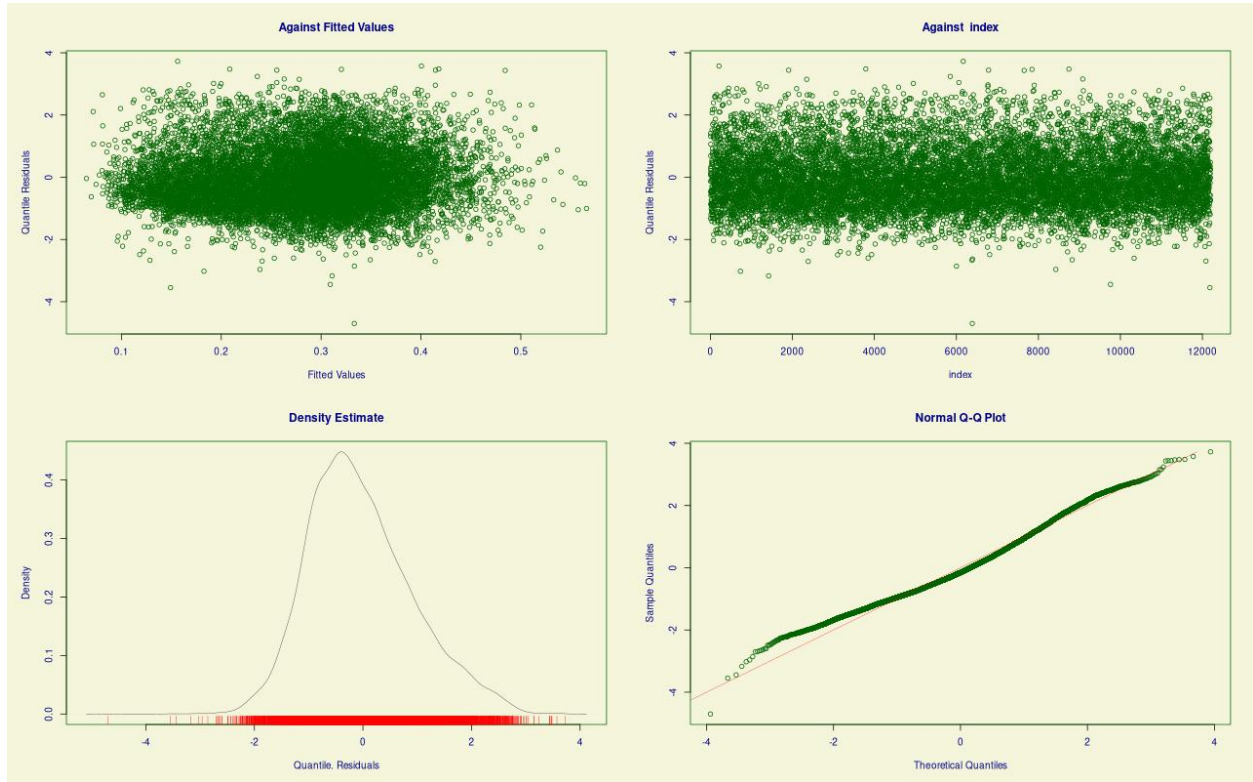


Figure 4: Diagnostics for beta inflated GAMLSS model predicting proportions of nominal -s

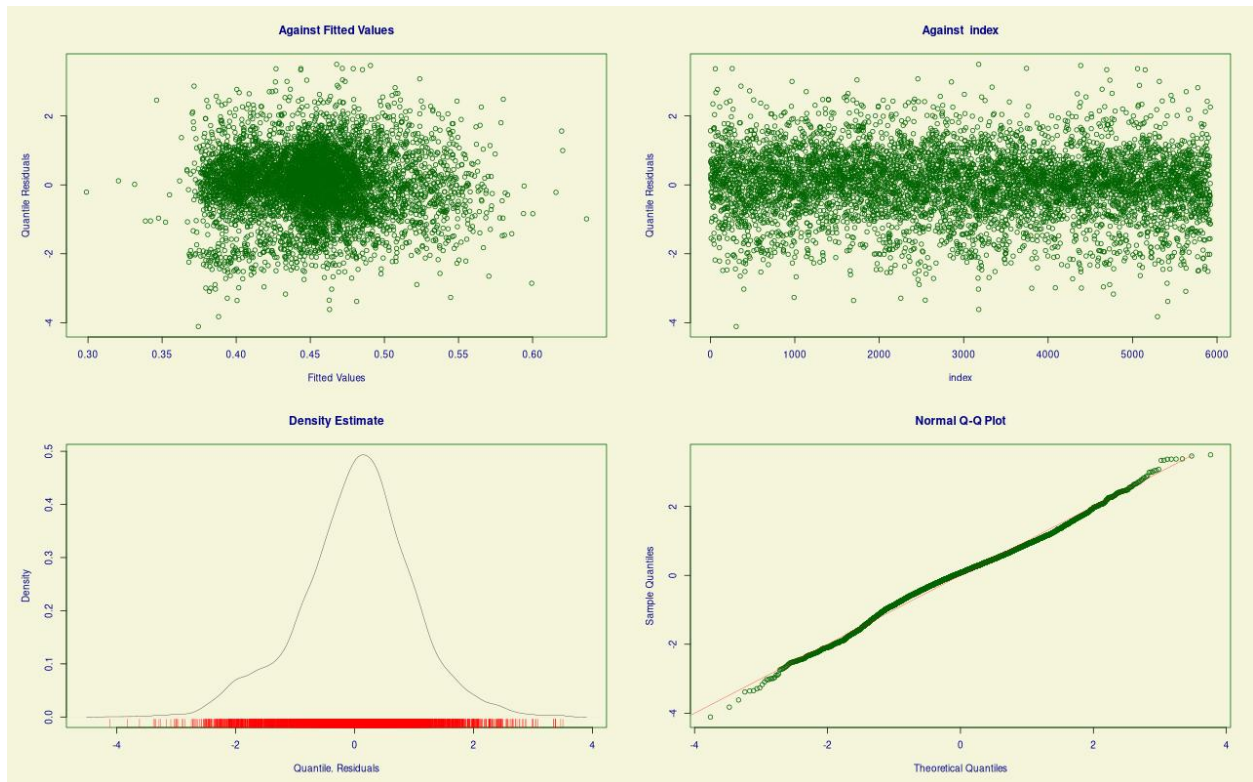


Figure 5: Diagnostics for beta inflated GAMLSS model predicting proportions of -ed

Bibliography

- Beekhuizen, Barend, Blair C. Armstrong & Suzanne Stevenson. 2021. Probing lexical ambiguity: Word vectors encode number and relatedness of senses. *Cognitive Science* 45(5). e12943. <https://doi.org/10.1111/cogs.12943>. <https://onlinelibrary.wiley.com/doi/abs/10.1111/cogs.12943>.
- Diessel, Holger. 2016. Frequency and lexical specificity in grammar: A critical review. In Heike Behrens & Stefan Pfänder (eds.), *Experience counts: Frequency effects in language*, 209–238. De Gruyter. <https://doi.org/10.1515/9783110346916-009>.
- Dowle, Matt & Arun Srinivasan. 2021. *Data.table: Extension of ‘data.frame’*. <https://CRAN.R-project.org/package=data.table>.
- Du, Jiaju, Fanchao Qi & Maosong Sun. 2019. Using BERT for word sense disambiguation. *CoRR* abs/1909.08358. <http://arxiv.org/abs/1909.08358>.
- Evert, Stefan. 2005. *The statistics of word cooccurrences: Word pairs and collocations*. Friedrich-Alexander-Universität Erlangen-Nürnberg PhD thesis. <http://www.collocations.de/phd.html>.
- Evert, Stefan & Andrew Hardie. 2011. Twenty-first century corpus workbench: Updating a query architecture for the new millennium. University of Birmingham.
- Gries, Stefan. 2021. A new approach to (key) keywords analysis: Using frequency, and now also dispersion. *Research in Corpus Linguistics* 9(2). 1–33. <https://doi.org/10.32714/ricl.09.02.02>.
- Gries, Stefan Th. 2008. Dispersions and adjusted frequencies in corpora. *International journal of corpus linguistics*. John Benjamins 13(4). 403–437.
- Kottur, Satwik, Ramakrishna Vedantam, José MF Moura & Devi Parikh. 2016. Visual word2vec (vis-w2v): Learning visually grounded word embeddings using abstract scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4985–4994.
- Kromer, Victor. 2003. A usage measure based on psychophysical relations. *Journal of Quantitative Linguistics*. Taylor & Francis 10(2). 177–186.
- Lee, Younghoon. 2021. Systematic homonym detection and replacement based on contextual word embedding. *Neural Processing Letters* 53(1). 17–36. <https://doi.org/10.1007/s11063-020-10376-8>.
- Lijffijt, Jeffrey & Stefan Th Gries. 2012. Correction to stefan th. Gries “dispersions and adjusted frequencies in corpora” international journal of corpus linguistics 13: 4 (2008), 403-437. CiteSeer.
- Love, Robbie, Claire Dembry, Andrew Hardie, Vaclav Brezina & Tony McEnery. 2017. The spoken BNC2014: Designing and building a spoken corpus of everyday conversations. *International Journal of Corpus Linguistics*. John Benjamins 22(3). 319–344.
- Ospina, Raydonal & Silvia L. P. Ferrari. 2012. A general class of zero-or-one inflated beta regression models. *Computational Statistics & Data Analysis* 56(6). 1609–1623. <https://doi.org/10.1016/j.csda.2011.10.005>. <https://www.sciencedirect.com/science/article/pii/S0167947311003628>.
- Pennington, Jeffrey, Richard Socher & Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 1532–1543.
- Piantadosi, Steven T., Harry Tily & Edward Gibson. 2012. The communicative function of ambiguity in language. *Cognition* 122(3). 280–291. <https://doi.org/10.1016/j.cognition.2011.10.004>. <https://www.sciencedirect.com/science/article/pii/S0010027711002496>.
- Plag, Ingo, Julia Homann & Gero Kunter. 2017. Homophony and morphology: The acoustics of word-final s in english. *Journal of Linguistics*. Cambridge University Press 53(1). 181–216.
- R Core Team. 2021. *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Rigby, R. A. & D. M. Stasinopoulos. 2005. Generalized additive models for location, scale and shape, (with discussion). *Applied Statistics* 54.3. 507–554.
- Rigby, Robert A, Mikis D Stasinopoulos, Gillian Z Heller & Fernanda De Bastiani. 2019. *Distributions for modeling location, scale, and shape: Using GAMLSS in r*. Chapman; Hall/CRC.
- Selivanov, Dmitriy, Manuel Bickel & Qing Wang. 2020. *text2vec: Modern text mining framework for r*. <https://CRAN.R-project.org/package=text2vec>.
- Shahmohammadi, Hassan, Hendrik Lensch & R Harald Baayen. 2021. Learning zero-shot multifaceted visually grounded word embeddings via multi-task training. *arXiv preprint arXiv:2104.07500*.
- Stasinopoulos, Mikis D, Robert A Rigby & Fernanda De Bastiani. 2018. GAMLSS: A distributional regression approach. *Statistical Modelling*. SAGE Publications Sage India: New Delhi, India 18(3-4). 248–273.

- Stefanowitsch, Anatol & Stefan Th Gries. 2003. Collostructions: Investigating the interaction of words and constructions. *International journal of corpus linguistics*. John Benjamins 8(2). 209–243.
- The British National Corpus, version 3 (BNC XML Edition). 2007. <http://www.natcorp.ox.ac.uk/>; Distributed by Bodleian Libraries, University of Oxford, on behalf of the BNC Consortium.
- Tomaschek, Fabian, Ingo Plag, Mirjam Ernestus & R. Harald Baayen. 2021. Phonetic effects of morphology and context: Modeling the duration of word-final s in english with naïve discriminative learning. *Journal of Linguistics*. Cambridge University Press 57(1). 123–161. <https://doi.org/10.1017/S0022226719000203>.
- Trott, Sean & Benjamin Bergen. 2020. Why do human languages have homophones? *Cognition* 205. 104449. <https://doi.org/10.1016/j.cognition.2020.104449>. <https://www.sciencedirect.com/science/article/pii/S0010027720302687>.
- Wasow, Thomas. 2015. Ambiguity avoidance is overrated. In Susanne Winkler (ed.), *Ambiguity*, 29–48. De Gruyter. <https://doi.org/10.1515/9783110403589-003>.
- Wickham, Hadley. 2016. *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wiedemann, Gregor, Steffen Remus, Avi Chawla & Chris Biemann. 2019. Does BERT make any sense? Interpretable word sense disambiguation with contextualized embeddings. *arXiv preprint arXiv:1909.10430*.
- Wood, S. N. 2017. *Generalized additive models: An introduction with r*. 2nd edn. Chapman; Hall/CRC.
- Wood, S. N., N., Pya & B. S'afken. 2016. Smoothing parameter and model selection for general smooth models (with discussion). *Journal of the American Statistical Association* 111. 1548–1575.
- Yung Song, Jae, Katherine Demuth, Karen Evans & Stefanie Shattuck-Hufnagel. 2013. Durational cues to fricative codas in 2-year-olds' american english: Voicing and morphemic factors. *The Journal of the Acoustical Society of America*. Acoustical Society of America 133(5). 2931–2946.
- Zimmermann, Richard. 2020. Word growth dispersion — a single corpus part measure of lexical dispersion. Paper presented at ICAME41 Heidelberg. <https://www.youtube.com/watch?v=k8etOvRcF4c>.