The University of
Nottingham

UNITED KINGDOM · CHINA · MALAYSIA

Smith, Catherine and Adolphs, Svenja and Harvey, Kevin and Mullany, Louise (2014) Spelling errors and keywords in born-digital data: a case study using the Teenage Health Freak Corpus. Corpora, 9 (2). pp. 137-154. ISSN 1755-1676

**Access from the University of Nottingham repository:**
http://eprints.nottingham.ac.uk/35782/1/Smith%2C%20Adolphs%2C%20Harvey%2C%20and%20Mullany_Spelling%20Errors%20and%20Keywords%20in%20Born-Digital%20Data.pdf

Spelling Errors and Keywords in Born-Digital Data: A Case Study using the Teenage Health Freak Corpus

**Abstract**
The abundance of language data now available in digital form and the rise of particular language varieties used for digital communication means that issues of non-standard spelling and spelling errors are likely to become a more prominent issue for compilers of such corpora. This paper examines the effect of spelling variation on keywords in a born-digital corpus in order to explore the extent and impact of this variation for future corpus studies. The corpus used in this study consists of emails about heath concerns sent to a health website by adolescents. Keywords are generated using the original version of the corpus and a version with spelling errors corrected with the BNC as the reference corpus. The ranks of the keywords are shown to be very similar and therefore suggest that, depending on the research goals, keywords could be generated reliably without any need for spelling correction.

## 1. Introduction

Corpora are more frequently being created from texts taken from the web and therefore the issue of spelling errors and non-standard spelling are an increasing issue for corpus compilers. Unlike corpora built with published, well-edited materials the increasing spread of user generated content present in the web 2.0 world presents new challenges to the corpus creator in terms of *typos* as this material is not edited to the standard of traditional print media or even large company websites. In addition the same way that non-standard spelling poses problems for corpus researchers working with historical corpora (see Baron et al. 2009) the use of non-standard spelling in e-language is reintroducing the problem for modern corpora. The use of chat and text abbreviations is widespread in the internet community and innovations are always being introduced (Crystal 2006, 2011; Baron 2008). This paper aims to investigate the impact that spelling errors and non-standard spelling may have on keywords generated with a born-digital corpus.

The corpus used in this case study is comprised of health questions sent to Doctor Ann through the *Ask Doctor Ann* facility on the Teenage Health Freak website.[1] The messages date from January 2004 to December 2009, a period of six years. With the exception of the removal of very similar messages sent within a small time frame the corpus is unedited and contains all messages sent to the website during the time period in question, only a small number of which are published on the website with answers from the medical staff. The messages themselves are typed directly into a web form and other than the automatic removal of personal details such as email addresses the corpus contains exactly what was typed into the form by the users. In total the corpus contains 113,480 messages and

[1] http://www.teenagehealthfreak.org

2,217,919 words. This corpus is unique in many respects and is of particular interest to the medical community since the questions posed are unsolicited and the users of the site are able to ask about whatever area of their health is concerning them. The fact that the corpus is generated by adolescents is also of interest to the linguistic community. Adolescents are typically seen as language innovators (Stenström et al. 2002) and this is no exception in the on-line community (Crystal 2006:94)

Keywords are the starting point for a great many lexically oriented corpus studies (Scott and Tribble, 2006: 55-72; Baker 2006; Harvey et al. 2008). They provide a good segue into the study of a large corpus highlighting words and topics which are unusually frequent in the data and therefore may warrant further investigation. Such studies are classed by Rayson as type III corpus studies which are categorised by 'the use of corpus-based comparative frequency evidence to drive the selection of words for further study' (Rayson, 2008: 523). Scott defines a keyword as 'a word which occurs with unusual frequency in a given text... by comparison with a reference corpus of some kind' (Scott, 1997: 236). Statistical procedures are used to determine whether the comparative frequencies of a word in the target and reference corpora are significant enough to classify the word as key in the target corpus. Several measures have been used to determine which words should be considered key including the chi-squared test (Hofland and Johansson, 1982), Mann-Whitney test (Kilgarriff, 2001) and the T-score (Paquot and Bestgen, 2009). The most widely used statistical procedure however is log-likelihood.

The log-likelihood statistic is calculated using the total counts for the word in question in each of the corpora, the observed frequency, and the expected frequency which is calculated to take account of the size of each of the corpora. Fundamental to the log-likelihood statistic, as with all other measures mentioned above, are the word counts for each of the words in both the target and reference corpora. If one of the corpora has problems with non-standard spelling or contains a large volume of spelling errors the counts for the words will be affected and consequently also the log-likelihood score. If for example the word *what* occurs in the corpus with the standard spelling and also with the two chat style abbreviations *wat* and *wot* then the overall counts for *what* in the target corpus will be reduced as it is being represented by several different orthographical forms. This problem becomes even greater for the statistical calculation if, as in the case here, only the target corpus contains such inconsistencies. In this case the non-standard and incorrect spellings could dominate the keyword list as they are not likely to be present in the reference corpus at all. This kind of challenge is encountered in a variety of different corpora: historical corpora, for example, in which spelling conventions were yet to be established (Baron et al. 2009), also regional corpora where dialects are being represented (Kay, 2006) as well as born digital corpora which use non-standard spelling as used in this research.

This paper first analyses the volume and type of spelling errors and non-standard spellings found in the Teenage Health Freak corpus to establish the scale and nature of the problem faced. The variant spelling in the corpus is then corrected as far as possible to create a second normalised corpus. These two corpora are then compared with respect to the keywords they produce against the same reference corpus in order to establish any differences in keyword rank due to the spelling variation. In this paper the term spelling error is used as an umbrella term for spelling error, deliberate non-standard spelling and other abbreviations or acronyms used for words or phrases. This is done as a matter of convenience and is not intended as any judgement on the language use itself.

## 2. Corpus Approaches to Spelling Variants and Errors

Interest in spelling variation in corpus linguistics has typically been focussed on historical corpora and learner corpora. In historical corpora the problem is caused by a lack of standardised spelling. Research by Baron et al. suggests that variant tokens make up between 35 percent and 40 percent of all tokens in English corpora from the period 1400 to 1550 gradually reducing to below 10 percent by 1650 (2009: 53). The problems that this level of variation causes for part of speech tagging and subsequent semantic tagging using the online tool Wmatrix (Rayson 2009)[2] are what inspired the creation of a VARD, a variant detection tool (Archer et al., 2003; Rayson et al., 2007: 1). The latest version of VARD achieved impressive results with historical corpora with precision over 90 percent and recall reaching 65 percent with sufficient training (Baron and Rayson, 2009: 14). Recently VARD has been developed from a tool specifically designed to deal with historic English corpora to a tool that can potentially be used to deal with any kind of spelling variation or error in corpora written in any language (Baron and Rayson, 2009: 4-9, 13). VARD can be customised in several ways, for example by changing the dictionary of accepted spellings or adding new letter replacement rules, and can also be trained with a specific dataset (Baron and Rayson, 2009:9). The performance of VARD on a corpus of Children's writing was much lower than the figures achieved with historical corpora with a precision of around 80 percent but a recall no higher than 20 percent. In this experiment the only customisation performed was to train VARD with samples of the corpus and much better recall figures could be expected if all of the customisation features of VARD had been employed (Baron and Rayson, 2009:19).

In learner corpora spelling errors are only one of several types of errors that are of particular interest to the compilers and users of the data. Most of the work in identifying the errors is done manually but a few tools have also been developed to assist with the process of finding and correcting errors. A tool has been developed for computer assisted error annotation, UCLEE (Université Catholique de Louvain Error Editor), but it functions to speed up the manual process of mark-up rather than assist with the identification of errors themselves (Dagneaux et al., 1998 :167-168). Rayson and Baron have also tested VARD's performance on learner corpora to see if it could be used to aid with error annotation (2011). With training, as with the test on Children's writing, they were able to achieve a very high precision figure of 90.8 but the recall level was again much lower at 23.4. This is attributed in this case to the large number of real world errors found in learner corpora which are not yet handled by VARD (Rayson and Baron, 2011: 122).

Studies of modern born-digital data have tended not to address the issue of normalising spelling variation as the studies have typically focussed on the innovations of language and orthography used (Tagg, 2009; Ooi et al., 2007; Hoffman, 2007). In a study of Singaporean English blogs Ooi et al. find the same kind of reduction in performance when using Wmatrix's semantic tagger as was found with historical corpora. In this study Wmatrix is used to analyse two corpora one containing blogs written by teenagers and the other by undergraduates. In total 3,712 types from the undergraduate corpus were left unclassified by the semantic tagger and a much higher 11,137 types from the teenagers corpus. Rather than considering any kind of pre-processing to normalise the corpus Ooi et al. (2007) suggest that corpus tools like Wmatrix need to evolve to be able to deal with this new emerging form of English. Tagg (2009) is also interested in the creative use of language but in this case with a

---

[2] WMatrix is an online tool for the analysis and comparison of corpora. In particular it facilitates part-of-speech and semantic annotation of corpora which can then be analysed using a the same statistical procedure as is used for keywords.

corpus of text messages. Although Tagg chose not to normalise the language for this particular study she does note that some experimentation with retraining VARD on this data suggests that its use could be possible with the corpus. In a later conference presentation the latest version of VARD was used on the corpus with successful results (Tag et al., 2010).

**3. Spelling Errors in the Teenage Health Freak Corpus**
In order to gain an insight into the scale and type of spelling errors present in our corpus of health questions from teenagers, 50 messages were sampled at random from each year in the corpus (a total of 300 messages). These messages were manually checked for incorrect and unconventionally spelled words and the nature of each spelling error was analysed. The volume of spelling errors is summarised in Table 1. As the table shows based on this sample study it is estimated that 7.6 percent of the words in the Teenage Health Freak corpus could be incorrectly or unconventionally spelled. If the samples are representative of the whole corpus this would amount to around 168,600 words. It is not common to report data regarding spelling variation in electronic communication corpora but Tag et al. report figures from a text message corpora which are similar to those found in the Teenage Health Freak data. As part of a test of automated spelling normalisation with VARD a test sample is manually corrected. The sample consisted of 2,430 messages containing a total of 41,342 words. Of these words 3,166 required standardisation meaning 7.7 percent of the words in the sample were incorrectly or unconventionally spelled (Tagg et al., 2010). In addition both corpora found that while some messages contained a lot of unconventionally spelled words other messages had no variation at all.

Table 1

The errors can be classified into five broad categories: chat-style abbreviations; phonetic errors; typographical errors; deliberate errors for emphasis; and finally errors that didn't fit into any of the other categories. Chat-Style language includes abbreviations and acronyms which might be expected in text messages or instant messaging: in our samples these include the following transformations: u > you; 4 > for; cuz > because; sum > some. Typographical errors are errors which are most likely to be caused by mistyping and include the following examples: typographical iam > i am; resulst > results; alchohl > alcohol. Phonetic errors are words which can reasonably be pronounced in the same way as the original word and are less likely to have been caused by mistyping. In our sample this includes: probarbly > probably; egsisting > existing; marige > marriage. The next category of emphasis is less of an error and more of a manipulation of the language. However since this use of language poses the same problems for keyword analysis as the other spelling errors it is considered here. The emphasis category includes deliberate errors made for emphatic purposes. These are typically additions of letters as in these examples from our sample messages: soooo > so; yoooo > yo. The final category includes everything that doesn't fit into any of the other categories, in our samples this includes only one example:  pencise > penis.

Table 2

Table 2 shows the number of errors in each category present in our sample messages. As the table shows the vast majority of errors are accounted for by typographical errors and chat-style errors. The nature of chat-style errors should make correcting them relatively easy as there is a great deal of internal consistency with this type of language use (see Tagg, 2009:136-8 for a summary of the use of such language in a text message corpus). An

interesting observation on spelling in general but which is particularly true of the use of chat-style abbreviations is that there is huge difference between messages with some users avoiding all chat-style language and others making full use of it. This may reflect the familiarity of the user with instant messaging, forum writing and perhaps text messaging but also reflects the choice of register considered appropriate for addressing medical questions to Dr Ann which some selecting very formal registers and other much more informal. Compared to chat-style errors the typographical errors are not as consistent. However along with the phonetic errors the vast majority consist of a single omission, addition, deletion, substitution or transposition. This means that distance based spelling detection algorithms may do a reasonable job at identifying these errors (Jurafsky and Martin, 2000:144-6).

## 4. Spelling Correction Procedure

In order to try and correct the spellings in the corpus an evaluation of VARD was conducted as this was the program used to regularise the spelling in Baron et al. (2009). At the time the project started the latest version was VARD2.2 which was still specifically aimed at the regularisation of historical corpora although it had also been used for other language varieties. Communication with Baron suggested that VARD2.3 would be better suited to our data as many changes were being made in the automated processing to widen its application. During the course of the project VARD2.3 was released but unfortunately the release date was too far into the project for it to be used to correct the spelling for the research reported here.

Instead of using a specifically designed tool the spelling was corrected with the help of WordSmith Tools' keyword procedure (Scott 2008). The main advantage of this procedure is that it is well known among corpus linguists and is available in most of the standard software tools available and would therefore be easily replicable for other corpora. The written component of the BNC was selected as the reference corpus as this minimises the number of non-standard tokens present in the reference corpus. As described above, the keyword procedure works to highlight words which are unusually frequent in the target corpus in comparison to their frequency in the reference corpus. Therefore any words which are incorrectly spelled but occur in the same variant form frequently in our corpus were present in the keyword list. The resulting keyword list was reviewed manually and any variant forms which could be corrected to the same word in the vast majority of instances were corrected. This analysis was supported by the use of Concordance lines to allow the intended word to be determined from the context. In cases where the error could not, in most cases, be corrected to a single word the word was not corrected and the error remained. Once the whole list had been processed in this way a script was used to mark the errors and provide their corrections. As the corpus was in XML the spelling corrections were indicated using the Text Encoding Initiative (TEI) *choice*, *sic* and *corr* tags so that the original spelling was also preserved. The resulting XML for each corrected word looks like this:

```
<choice>
        <sic>plz</sic>
        <corr>please</corr>
</choice>
```

Using this method 2,732 types were corrected which amounts to 88,542 tokens, just over 50 percent of the predicted error total. The 2,732 incorrect types were corrected to just 900 types suggesting there could be a large impact on word frequencies. As missing

apostrophes have an effect on word tokenisation and therefore on frequencies and keywords, they were also counted as spelling errors in this study. 46 of the corrected types only involved missing apostrophes making the total number of types corrected that did not involve apostrophes 2,686. This has a large effect on the total tokens corrected removing 35,077 leaving only 53,465 corrected tokens. Because of the large number of changes just involving apostrophes the effects on keywords will be measured with and without the apostrophe corrections made.

## 5. The Effects of Spelling Errors on Key Words
### 5.1. Methodology
The procedure for comparing the keyword lists was based on that used by Baron et al. (2009) when comparing the effects of non-standardised spelling in early modern English on the results of keyword analysis. In this paper Baron et al. use two statistical measures, Spearman's rank correlation coefficient and Kendall's Tau rank correlation coefficient to compare two keyword lists. These statistical measures both focus on differences in ranks but they are sensitive to different types of movement within the ranked lists being compared. As Spearman's rho uses the squared difference in the ranks between the two lists (Conover, 1999: 287) this measure is more sensitive to movements over a greater distance within the ranked lists and relatively insensitive to large numbers of small movements within the ranks. In contrast Kenall's tau is based on concordant and discordant pairs (Sprent and Smeeton, 2007: 318) and is therefore more sensitive to the volume of changes rather than the actual difference between the two ranks and therefore the distance of the movement. In most cases Spearman has been found to give a slightly higher figure than Kendall (Sprent and Smeeton, 2007: 323)

It should be noted here that Scott in his "frequently asked questions" for WordSmith 5.0 advises that it is unsafe to rely on rank order in keyword lists (Scott, no date). As log-likelihood is a statistical measure of significance it is true that words are either key at the specified p value or are not key. However, in practice, researchers must often resort to relying on the top N keywords and in this regard changes in rank caused by spelling errors or a change in reference corpus could have an impact on the focus of the research. Changes in rank order also give a broad overview of the impact of such changes on the keyword lists.

The procedure used is outlined below.

- Generate keyword lists to compare using WordSmith Tools and specified parameters for corrected and uncorrected spelling
- Remove any words not present in both of the resulting keyword lists
- Rank remaining entries in both lists from 1 to n (n will be the same for both lists)
- Use the ranks as the input to correlation graphs and statistical procedures

### 5.2. Effect on keyword lists
The British National Corpus was chosen as the reference corpus to test the effect of spelling variants on the keyword lists. The keywords used in this section were all generated using WordSmith Tools (Scott, 2008) using the log-likelihood statistical measure. The minimum frequency threshold was 5 and the p value used was 0.000001. Keywords were generated against the reference corpus for the Teenage Health Freak corpus based on the original spelling, the fully corrected spelling, and the corrected spelling ignoring missing apostrophes. The first thing to look at is the number of keywords generated for each version of the corpus which is shown in Table 3 below.

Table 3

By correcting the spelling in the corpus the number of keywords generated is reduced by over 1,500. An examination of the words that are key only when the spelling errors have been corrected show that they contain medical or medical related terms. In all of these cases however, they occur in some alternatively spelled form in the keyword list based on the uncorrected spelling so are not lost entirely even when the spelling is not corrected. Also as Table 4 shows these words occur a long way down the keyword lists for both the original and corrected spelling.

Table 4

The other words that are only present in the corrected spelling list also appear a long way down the keyword list. Only four examples occur in the top 1,000 keywords and these also occur in various spellings in the keyword list generated from the original spelling. These words are shown in Table 5. In all these cases correcting the spelling from what is typically multiple incorrect forms when added together means that the frequency becomes high enough for the correctly spelled version to be considered key. Only *deodorant* occurs with the correct spelling in both keyword lists but the correctly spelled version occurs considerably lower in the original spelling list than the incorrect spellings do.

Table 5

The keywords which occur only in the list generated with original spelling are predominantly examples of non-standard spelling and with the exception of some of the standard chat-style abbreviations they tend to occur towards the bottom end of the keyword list. There would however be a lot more keywords to analyse had the spelling not been corrected. In some cases the correction of the spelling also makes a huge difference to the frequency count for the word and therefore its rank in the keyword list. An example of such a word is embarrassed which has multiple spellings as can be seen in Table 6.

Table 6

In the version of the corpus with the spelling corrected the word frequency for embarrassed is 596 making it 171 in the keyword list rank with a G2 value of 2,206.33. This is higher than the highest entry in the original spelling keyword list where *embarrased* is ranked 368 with a G2 value of 947.95 and considerably higher than the correctly spelled version which is ranked 1,134 with a G2 value of only 156.97. It seems then, that even if correctly spelled keywords are not lost altogether when the spelling errors remain in the corpus, the spelling variation could make a difference to the ranks of the keywords. This was the hypothesis which formed the basis of Baron et al.'s investigation into spelling variants in early modern English and will also form the basis of our investigation.

The Spearman rank correlation coefficient and the Kendall correlation coefficient were calculated between two pairs of keyword lists all generated using the BNC as the reference corpus. The pairs were the original spelling and the fully corrected spelling, and the

original spelling against the corrected spelling ignoring missing apostrophes.[3] An indication of the type of correlation present in the data can be seen in Figure 1 and Figure 2. These graphs suggest a very strong positive correlation in the keyword ranks. Anything above the perfect line of correlation represent words that are spelled incorrectly, but in the same way, in the Teenage Health Freak corpus so many times that the incorrect spelling is much higher in the uncorrected keyword list than it is when it is corrected and compared to the correctly spelled equivalent in the reference corpus. The strongest example of this in the corpus is the word *conscious* (point 473, 1209) which is misspelled as *concious* so frequently in the corpus that it ranks very high in the keyword list until it is corrected and therefore compared with the actual frequency of the word in the reference corpus. Anything below the perfect line of correlation represent words where, when all the spelling variants are corrected, they collectively make a big enough difference to the frequency for them to move higher up the keyword list for the corrected text, as with the example of *embarrassed*. It can also be seen from the graphs that the later items towards the bottom of the ranked list are more affected than those towards the top.

Figure 1

Figure 2

Table 7

Both Spearman's rho and Kendall's Tau return a number between -1 (a perfect negative correlation) and +1 (a perfect positive correlation) with 0 being no correlation at all. The results for these comparisons can be seen in Table 4. These numbers are both very close to +1 suggesting a very positive correlation between the ranks of the keyword lists. Not correcting missing apostrophes leads to a slightly better correlation than the fully corrected text but both are very high. The critical tables for Spearman's *r* stop at 30 degrees of freedom and since we have many more than 30 samples the result of Spearman's *r* can be converted to a *t*-value and that checked with N-2 degrees of freedom against the critical values for *t*. In both cases our degrees of freedom (1849 and 1893) are off the scale for p values of *t* value but both scores are much greater than the values needed at the highest degrees of freedom available and therefore we can conclude that there is a significantly high correlation to suggest that the spelling variants in the corpus make little difference to the ranks of the keyword lists. These figures are certainly higher than those reported by Baron et al. whose overall Spearman's rho score for the Innsbruck Letter corpus in manually standardised and original form was 0.705 and the Kendall's Tau 0.530. This is to be expected given the much higher rates of variation present in the Innsbruck Letter corpus. When comparing automatically standardised (using VARD) and original texts over several decades the both correlation scores were much higher and above 0.9 by the 1600s, figures are only reported in

---

[3] Butler (1985) says Spearman's *r* should not be used if there are "a large number of tied ranks" (Butler, 1985: 147). In place of Spearman's *r* the ranks should be used as input to the Pearson product-moment correlation coefficient (the result of this calculation is however still known as Spearman's *r*). In both cases here the number of tied ranks is less than 25 percent of the total pairs in the study with no indication of what a "large number" might be it was decided to calculate the statistic using both methods for completeness. In our data there proved to be no difference in the calculations until the 7[th] or 8[th] decimal place so it was decided to report the results from the straight Spearman's *r* calculation for simplicity.

graphical form in this part of the study so it is not possible to give the highest correlation reached in that study however Spearman is very close to +1 and Kendall over 0.9 (Baron et al. 2009:58).

## 6. Conclusion

This case study demonstrates that even with born-digital data where the instances of spelling errors and non-standard spelling are likely to be higher than in other written corpora (i.e. newspaper or other print media corpora) spelling variation does not necessarily have a large impact on the keywords generated. Although the volume of keywords generated was much greater with the uncorrected version of the corpus the differences in the ranks of shared words were very small. The words found to be key in the uncorrected keyword list and not in the corrected keyword list tended to occur towards the bottom. They would, therefore, less of a problem for research focussing on words towards the top of the list. In general we can conclude that while correcting the spelling made a large difference to the ranks of some words (for example *embarrassed* noted above) the overall effect on the keyword ranks is small. So together with the fact even the correctly spelled technical terminology only present in the corrected keyword list still appears in some form on the uncorrected list, this suggests that for born-digital data the correction of spelling is not necessarily a central task before corpus analysis can take place. Of course for other types of corpus analysis on born digital data, such part-of-speech and semantic tagging that prompted the creation of a tool like VARD, where spelling variation is a much more significant problem, might suffer the same loss of accuracy as has been found with historical corpora.

More work is needed to establish whether the figure of 7% for non-standard spelling in the Teenage Health Freak corpus and in Tagg's text message corpus holds true for other CMC corpora. This will be a key factor as to how generalisable the results of this study are. While the range of messages in Teenage Health Freak corpus is quite wide there is a focus on health concerns which is unlikely to be found in more general CMC corpora and this topical focus may also contribute to the results seen in this study. The same research should be repeated on more general CMC corpora to support the suggestion that spelling correction is not of central importance in keyword studies of such data. The effort involved in such an investigation is somewhat lessened now in the era of internet corpora and 'big data'. In addition it is possible that a more accurate method of error detection and correction might show a much greater variation in key word rank since in this study the level of correction was around 50 percent. Now that the VARD has been adapted to deal well with spelling variation of non-standard Modern English, specifically CMC, it is possible that repeating this study on a version of the corpus that has been automatically corrected using VARD would shed more light on the impact of spelling on keyword rank. However, overall we suggest that this study demonstrates that while non-standard spelling should be a concern to corpus researchers working with born digital data, they can be confident in carrying out initial keyword based investigations on the uncorrected data.

References

Archer, D., T. McEnery, P. Rayson and A. Hardie. 2003. "Developing an automated semantic analysis system for Early Modern English". Archer, D., Rayson, P., Wilson, A. and T. McEnery (eds) *Proceedings of the Corpus Linguistics 2003 conference*. University Centre for Computer Corpus Research on Language: Lancaster University, pp. 22-31. Available at: http://ucrel.lancs.ac.uk/publications/CL2003/papers/archer.pdf (accessed November 2012).

Baker, P. 2006. *'The question is, how cruel is it?' Keywords, Foxhunting and the House of Commons*. Word Frequency and Keyword Extraction, AHRC ICT Methods Network Expert Seminar on Linguistics. Lancaster University. Available at http://www.arts-humanities.net/system/files/es1_07baker.pdf (accessed October 2011)

Baron, A. and P. Rayson. 2009. "Automatic standardization of texts containing spelling variation, how much training data do you need?". Mahlberg, M., V. González-Díaz, & C. Smith (Eds.), *Proceedings of the Corpus Linguistics Conference 2009*. Available at: http://ucrel.lancs.ac.uk/publications/cl2009/314_FullPaper.pdf (accessed October 2011).

Baron, N. S. 2008. *Always On: Language in an Online and Mobile World*. Oxford: Oxford University Press.

Baron, A., P. Rayson, & D. Archer. 2009. "Word Frequency and Keyword Statistics in Historical Corpus Linguistics". *Anglistik: International Journal of English Studies*, 20(1), 41-67.

Butler, C. 1985. *Statistics in Linguistics*. Oxford: Blackwell.

Conover, W.J. 1999. *Practical Nonparametric Statistics*. 3rd edn. New York: John Wiley & Sons.

Crystal, D. 2006. *Language and the Internet.* Cambridge: Cambridge University Press. 2nd edn.

Crystal, D. 2011. *Internet Linguistics: A Student Guide*. Oxford: Routledge.

Dagneaux, E., S. Denness, & S. Granger. 1998. "Computer-aided error analysis" *System,* 26, 163-174.

Harvey K, Churchill D, Crawford P, Brown B, Mullany L, Macfarlane A and McPherson A. 2008. Health communication and adolescents: what do their emails tell us? *Family Practice*. 25(4), 304-311.

Hoffman, S. 2007. "Processing Internet-derived Text - Creating a Corpus of Usenet Messages". *Literary and Linguistic Computing,* 22(2), 163-174.

Hofland and Johansson, 1982 *Word frequencies in British and American English*. Bergen: The Norwegian Computing Centre for the Humanities.

Jurafsky, D. & J. H. Martin. 2000. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition (International Edition)*. New Jersey, Prentice Hall.

Kay, C. 2006. *Issues for Historical and Regional Corpora: First Catch Your Word*. Word Frequency and Keyword Extraction, AHRC ICT Methods Network Expert Seminar on Linguistics. Lancaster University. Available at http://www.arts-humanities.net/system/files/es1_04kay.pdf (accessed October 2011)

Kilgarriff, A. 2001. "Comparing Corpora". *International Journal of Corpus Linguistics* 6:1, pp. 1-37.

Ooi, V. B. Y., P. K. W. Tan, & A. K. L. Chiang. 2007. "Analyzing personal weblogs in Singapore English: the Wmatrix approach". *Studies in Variation, Contacts and Change in English,* 2. Available at: http://www.helsinki.fi/varieng/journal/volumes/02/ooi_et_al/ (accessed October 2011).

Paquot, M. and Y. Bestgen 2009. "Distinctive words in academic writing: a comparison of three statistical tests for keyword extraction". In A. Jucker, D. Schreier, M. Hundt, ed(s), *Corpora: Pragmatics and Discourse*, Amsterdam, Rodopi, p. 247-269. Available at http://sites.uclouvain.be/cecl/archives/PAQUOT_BESTGEN_2009_Distinctive_words_in_ac ademic_writing_ICAME2008.pdf (accessed October 2011).

Rayson, P. 2009. *Wmatrix: a web-based corpus processing environment*. Computing Department, Lancaster University. Available at: http://ucrel.lancs.ac.uk/wmatrix/ (accessed October 2011).

Rayson, P., D. Archer, A. Baron, & N. Smith. 2007. "Tagging historical corpora - the problem of spelling variation". *Proceedings of Digital Historical Corpora, Dagstuhl-Seminar 06491*. Available at:
http://www.comp.lancs.ac.uk/~paul/publications/rabs_extAbs_dagstuhl06.pdf (accessed October 2011).

Rayson, P. 2008. "From Key Words to Key Semantic Domains" *International Journal of Corpus Linguistics*, 13: 4. pp. 519-549.

Rayson, P. & A. Baron. 2011. "Automatic error tagging of spelling mistakes in learner corpora". In F. Meunier, S. De Cock, G. Gilquin, & M. Paquot (Eds.), *A Taste for Corpora. In honour of Sylviane Granger*, Studies in Corpus Linguistics, 45. John Benjamins: Amsterdam. 109-126.

Scott, M. 1997. "PC Analysis of Key Words -- and Key Key Words", *System*, Vol. 25, No. 1, pp. 1-13.

Scott, M. 2008. *WordSmith Tools version 5*. Liverpool: Lexical Analysis Software.

Scott, M. no date. *Wordsmith 5.0 Answers to FAQs*. Available at:
http://www.lexically.net/wordsmith/version5/faqs/answers.htm#different_keynesses (accessed October 2011).

Sprent, P. and N.C. Smeeton. 2007. *Applied Nonparametric Statistical Methods* 4<sup>th</sup> edn. London: Taylor & Francis.

Stenström, A., G. Andersen, I. K. Hasund. 2002. *Trend in Teenage Talk: Corpus Compilation, Analysis and Findings*. Amsterdam/Philadelphia: John Benjamins.

Tagg, C. 2009. *A Corpus Linguistics Study of SMS Text Messaging*. unpublished PhD thesis, University of Birmingham.

Tagg, C., A. Baron, P. Rayson. 2010. *"I didn't spel that wrong did i. Oops" Analysis and standardisation of SMS spelling variation*. ICAME 2010. Available at: http://comp.eprints.lancs.ac.uk/2310/1/CorTxt_and_VARD_-_ICAME_presentation-Final.pdf (accessed November 2012).

| Year | No. of words | No. of errors | Percentage errors |
|------|-------------:|--------------:|------------------:|
| 2004 | 1209 | 76 | 6.3 |
| 2005 | 1403 | 116 | 8.3 |
| 2006 | 758 | 89 | 11.7 |
| 2007 | 898 | 55 | 6.1 |
| 2008 | 1000 | 70 | 7.0 |
| 2009 | 871 | 60 | 6.9 |
| All Years | 6139 | 466 | 7.6 |

Table 1: volume of spelling errors in 300 sample messages from the THF corpus

| Error Class | Total Occurrences |
|---|---|
| Typographical | 257 (ignoring apostrophes 134) |
| Chat-Style | 125 |
| Phonetic | 83 |
| Emphasis | 3 |
| None | 1 |

Table 2: classification of spelling errors in 300 sample messages from the THF corpus

| | Original Spelling | Corrected (ignoring apostrophes) | Corrected Spelling |
|---|---|---|---|
| Total keywords generated against BNC | 3,608 | 1,934 | 1,900 |

Table 3: Total number of keywords generated

| Corrected Spelling | | Original Spelling | |
|---|---|---|---|
| Spelling | Rank in List | Spelling | Rank in List |
| Transsexual | 952 | Transexual | 1038 |
| Achy | 1492 | Achey | 3415 |
| Tonsillitis | 1498 | Tonsilitis | 3389 |
| Bingeing | 1623 | Binging | 2466 |
| Syphilis | 1658 | Syphillis | 1673 |
| Dizziness | 1609 | Dizzyness | 2885 |
| Oestrogen | 1731 | Estrogen | 2048 |
| Tetanus | 1778 | Tetnus | 2523 |
| Disease | 1850 | Disese/Diseas | 2885/3415 |
| Lymph | 1898 | Lymphnodes | 2523 |

Table 1: Medical terms from the keyword List based on the corrected corpus and their equivalents in the original corpus

| Corrected Spelling | | Original Spelling | |
|---|---|---|---|
| Spelling | Rank in List | Spelling | Rank in List |
| Deodorant | 788 | Deodrant/Deoderant/Deodorant | 1313/1784/3526 |
| Accidentally | 903 | Accidently/Accidentaly | 887/2885 |
| Regularly | 980 | Regulary/Regulaly | 672/2885 |
| Noticeable | 991 | Noticable/Noticible | 795/2279 |

Table 2: Words that in the top 1,000 keywords from the corrected corpus and their equivalents in the original corpus

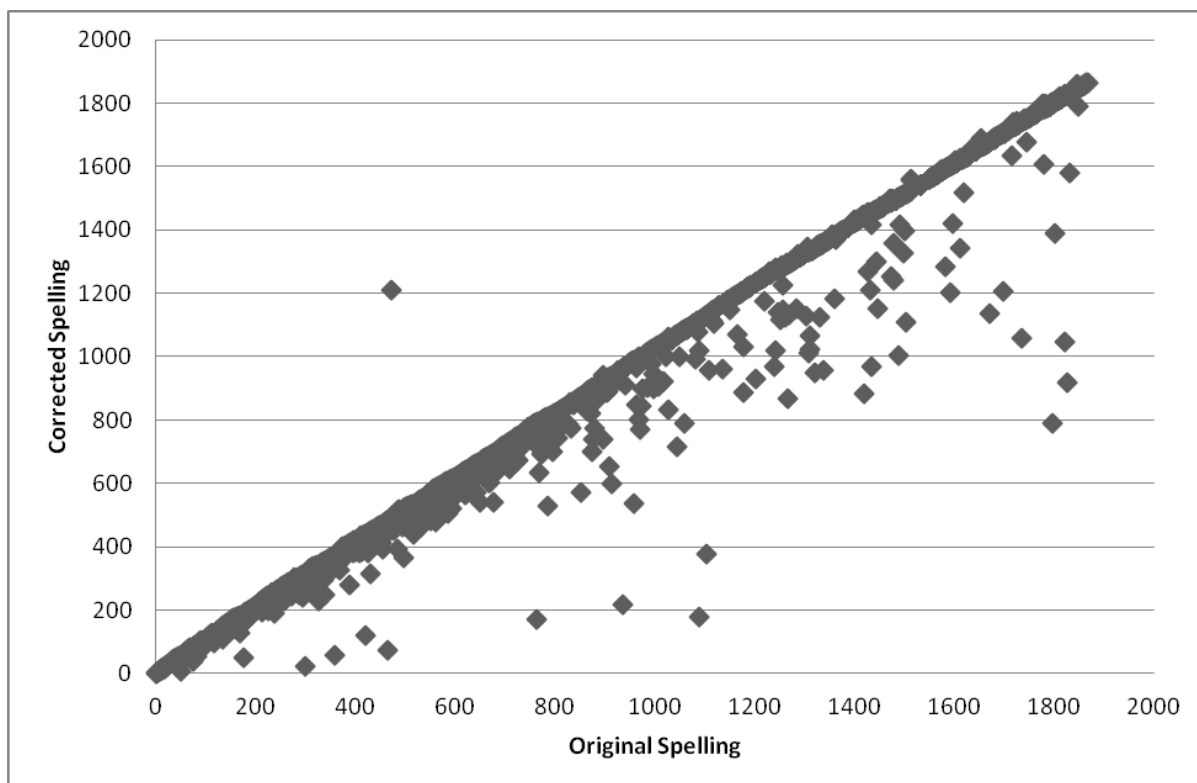| Spelling variant | Frequency | G2 | Rank |
|---|---|---|---|
| Embarrased | 125 | 947.95 | 368 |
| Embarassed | 108 | 801.61 | 413 |
| Embaressed | 63 | 483.62 | 578 |
| Embarresed | 60 | 460.59 | 589 |
| Embrassed | 21 | 161.21 | 1,104 |
| Embarrassed | 125 | 156.97 | 1,134 |
| Embarased | 14 | 107.47 | 1,416 |
| Embarrsed | 11 | 84.44 | 1,673 |
| Embaresed | 9 | 69.09 | 1,907 |
| Embarrised | 9 | 69.09 | 1,907 |
| Embarested | 5 | 38.38 | 2,885 |
| Embarised | 5 | 38.38 | 2,885 |
| Embarresd | 5 | 38.38 | 2,885 |
| Embarsed | 5 | 38.38 | 2,885 |
| Imbarrased | 5 | 38.38 | 2,885 |

Table 3: Variant spellings of embarrassed in the corpus

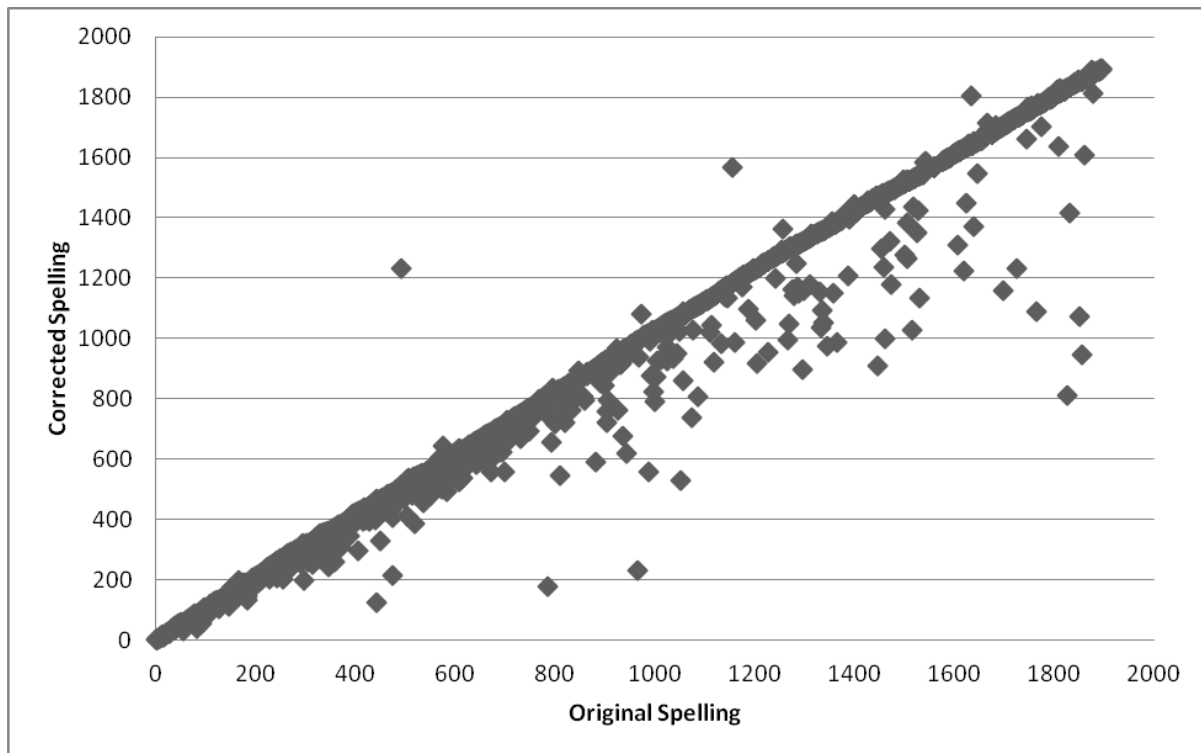Figure 1: correlation graph for fully corrected spelling against original spelling

Figure 2: correlation graph for corrected spelling ignoring missing apostrophes against original spelling

| | Fully Corrected Spelling against Original Spelling | Corrected Spelling Ignoring Apostrophes against Original Spelling |
|---|---|---|
| Kendall's Tau | 0.961 | 0.963 |
| Spearman's $r$ | 0.989 | 0.991 |
| Spearman's $r$ converted to t-value | 293 | 315 |

Table 4: Results from Spearman's rank correlation coefficient and its conversion to t-value