# Automatic Standardization of Spelling for Historical Text Mining

Alistair Baron, Paul Rayson (Lancaster University)
Dawn Archer (University of Central Lancashire)

## 1. Introduction

The use of textual data in humanities research is significantly aided by automated techniques such as key word analysis, collocations and corpus annotation (e.g. part-of-speech). If a text corpus contains a large amount of spelling variation, there is a considerable impact on the accuracy of these automatic techniques. For example, studies in respect to Early Modern English (EModE) corpora - the focus of the study detailed in this paper - have documented the adverse effects of spelling variation on key word analysis (Baron et al., forthcoming), part-of-speech tagging (Rayson et al., 2007) and semantic analysis (Archer et al., 2003).

The problem of spelling variation in corpora needs to be addressed in order for more accurate and meaningful results to be achieved in fields where historical source texts are required. Researchers can side-step the issue by using modernized versions of corpora, of course, but these are not always available. Another potential solution is to manually standardize the spellings; this includes reading through texts, spotting any non-standard spellings and deciding upon a modern equivalent, resulting in the production of a new version of the text with spelling variants replaced. However, a manual standardizing approach is likely to be unworkable when working with some of the larger corpora or online databases that are now available.

This paper details the current version of the VARiant Detector (VARD 2) tool, which can be used in various ways to standardize spelling variation in corpora. In particular,

the tool can be used to (partially) standardize spellings automatically, with no restriction on the number of words to be processed. Here, we focus on the ways in which the tool can be trained from manually standardized corpora samples, particularly the letter replacement component of the tool, and evaluate the improvement that this makes to the performance of VARD 2.

## 2. Early Modern English Spelling Variation

The EModE period is of particular interest in historical text mining studies; book production increased sharply during the period, largely due to the introduction of the printing press (1476) and increasing literacy levels (Görlach, 1991: 6). As such, the EModE period is the earliest period of the English Language from which a reasonably large corpus can be constructed and studied in detail.

Spelling variation was a major feature of EModE texts, the extent of which has recently been quantified in Baron et al. (forthcoming). It is common to find words spelt in a number of different forms in the same text or even on the same page. This was not seen as problematic, however, as there was no notion of the importance for a single spelling for each word; for example, letters would be added or removed to ease line justification. Vallins and Scragg (1965) and Culpeper and Archer (forthcoming) describe the spelling variation trends in more detail.

The effect of this spelling variation on textual analysis techniques has been shown in previous and forthcoming papers: key word analysis (Baron et al, forthcoming), part-of-speech tagging (Rayson et al., 2007) and semantic tagging (Archer et al., 2003). All of the studies showed that spelling variation causes considerable problems to the accuracy and meaningfulness of results, and that dealing with spelling variation (even partially) can achieve substantial improvements in annotation accuracy.[1] The

production of standardized or modernized versions of historical corpora therefore allows for more accurate automated text mining techniques to be applied to the corpora.[2]

## 3. VARD 2 and DICER

Our solution to the spelling variation problem described in the previous section has been the development of the VARD 2 tool,[3] a piece of software designed to assist researchers in standardizing historical corpora (specifically EModE texts) both manually and automatically. An earlier version of the VARD 2 software is described and evaluated in Rayson et al (2008). The current version can cater for user-created letter replacement rules, which will be used by the tool to find potential variant replacements. In addition, XML provision has been improved, processing speed increased, and a new word reference list[4] added. Screenshots of the latest version, VARD 2.2, are shown in Fig. 1 and Fig. 2.
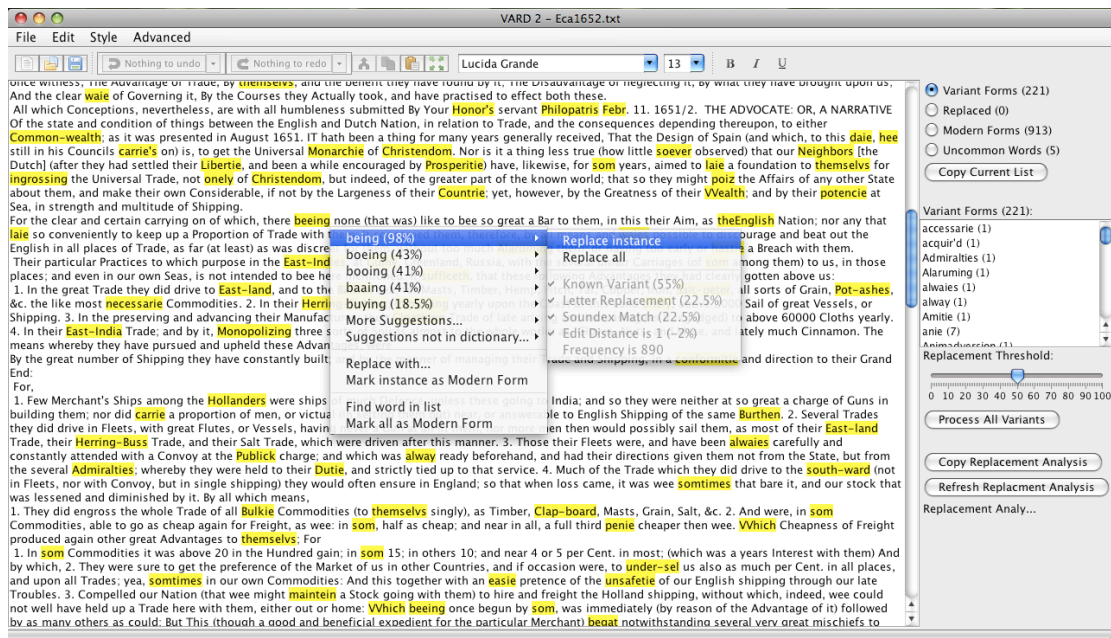


**Fig. 1 Screenshot of VARD 2.2 showing the interactive mode which allows the user to manually standardize texts and train the tool on samples of a corpus**
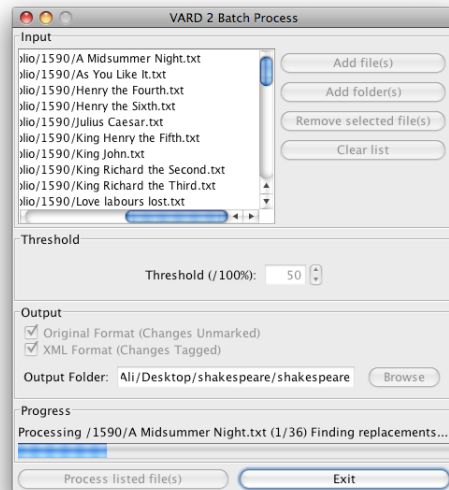
**Fig. 2 Screenshot of VARD 2.2 showing the batch-processing mode where users can automatically standardize a chosen group of texts**

The most useful way in which VARD 2 can be used is to automatically standardize spelling variation in an entire corpus. For EModE texts, this can be done immediately, with no training. However, for better results and to use the tool with other varieties of English, the user can train the software on a particular corpus by using the interactive version to manually process samples from the corpus. The tool will improve its ability to deal with a corpus based on decisions made by the user in the interactive version. It does this by learning which of its methods are most successful in finding the correct replacement for variants and adjusting its method weights accordingly (these are used when ranking potential replacements). The tool will also edit its dictionary and its list of specific variant replacements based on changes made by the user.

A new development to allow for further training of VARD 2 on a corpus is a tool named DICER (Discover and Investigation of Character Edit Rules). DICER can search XML output from VARD 2 for variant – replacement mappings or be provided with a list of such mappings. Each mapping is analyzed and a set of character edit rules are produced which can transform the spelling variant into its modern equivalent. The details of these character edit rules are then collated into a database,

which can be viewed through a set of web pages.[5] The main table produced by the analysis, shown in Fig. 3, displays details of the individual character edit rules along with various frequencies. By clicking on individual rules, further information is available such as which characters typically occur before and after the rule occurs; this is shown for the rule 'Delete e' in Fig. 4. Any frequency in the tables can be clicked to view a list of occurrences producing that frequency. The data produced in the DICER analysis is vast and thus cannot be detailed in full here.
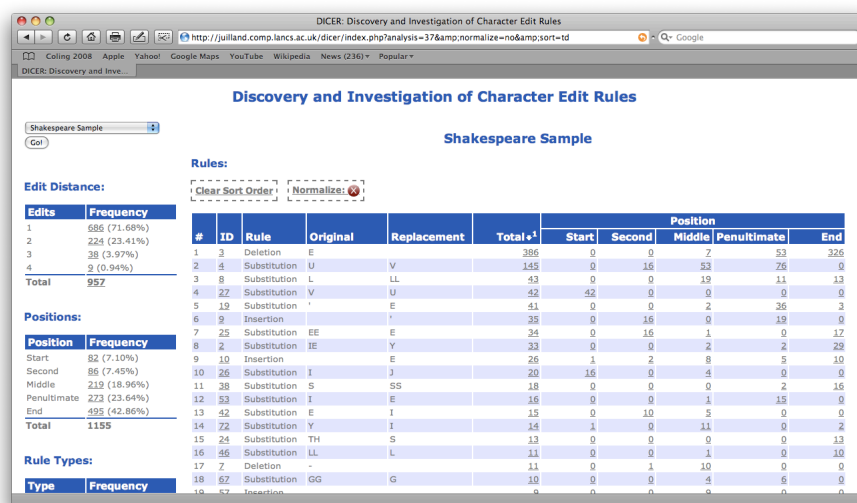


**Fig. 3 Screenshot of DICER analysis on a manually standardized 5,000-word sample of Shakespeare's First Folio**
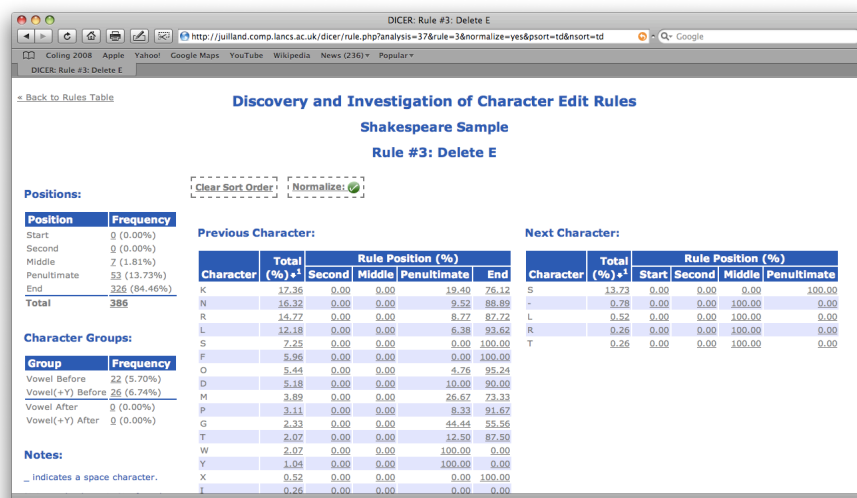


**Fig. 4 Screenshot of DICER showing the rule 'Delete E' in a manually standardized 5,000-word sample of Shakespeare's First Folio**

By using DICER to analyze manually standardized samples of a corpus, a list of common character edit rules can be viewed. These character edit rules can then be added to VARD 2 and the tool will be better equipped to make judgments on the correct replacement for variants found whilst automatically standardizing the corpus.

In order to test VARD 2 and DICER's training ability a 5,000-word sample of Shakespeare's First Folio[6] was manually standardized in the interactive-mode of VARD 2 as training data, the entire corpus was then automatically standardized. Using this small amount of training data (6% of the entire corpus) increased the proportion of spelling variants replaced[7] from 70.33% to 73.75%.

The automatic standardization (after training) resulted in 10,601 unique variant replacements. 70.35% of these replacements could be achieved through VARD 2's original set of character edit rules alone. DICER analysis was then produced on the manually standardized Shakespeare sample; this is shown in Fig. 3. VARD 2's rule list was augmented with additional rules from the DICER analysis: any rule occurring 10 or more times was added, if not already present. Using this new rule list 77.66% of the 10,601 unique replacements could now be found, an increase of 7.31%.

The results are extremely promising, and increasing the size of the manually standardized sample should improve these figures even further. DICER can also be used to provide probabilities dictating how likely a character edit rule should be applied in a certain position with specified surrounding characters. Modifying VARD to use these probabilities could see even greater improvements in performance.

## 4. Conclusion

This paper has described the problems that variant spellings cause for historical text mining, particularly for automated methods in historical corpus linguistics, such as part-of-speech tagging and key words analysis. In previous and forthcoming work, we have quantified the errors or differences that result from the application of untrained tools and techniques on historical data that has not been standardized. Our proposed solution is the VARD tool, which offers the potential to standardize spelling in historical texts automatically and with high accuracy. We have described recent improvements to VARD 2, such as the inclusion of a much larger modern dictionary that enables better detection of historical variants and matching with modern forms.

VARD 2 has been developed to deal with spelling variation in EModE texts; the tool can be used with its default settings to achieve partial standardization automatically. However, with some training, we have shown that VARD 2's performance is enhanced. Further training could allow the tool to be used with other varieties of non-standard English (e.g. SMS corpora and weblogs).

In the future, we will evaluate the extent to which variation that can only be detected contextually (e.g. 'then' for 'than' and 'bee' instead of 'be') contributes to the problem. Dealing with this problem requires more advanced techniques, e.g. POS tagging, to be used in the detection phase.

## Notes

1. Of course, spelling variants themselves are important linguistic features and thus worthy of study: as such, although our focus relates to how we might deal with spelling variation within historical data as a means of enabling the (more)

effective use of automated analytical techniques, we advocate that any solution to this problem should always retain the original spelling.

2. It should be noted that the accuracy of annotation is likely to be affected by additional factors, including differences in the grammar of the EmodE period when compared to present-day English (see Kytö and Voutilainen, 1995) and the possibility of a semantic shift in words from EModE to present-day English (see, for example, Knapp, 2000).

3. The tool is available to download online, with a user guide also provided. The software is free to use for academic purposes from http://www.comp.lancs.ac.uk/~barona/vard2/

4. Derived from the Spell Checking Oriented Word List (SCOWL). See http://wordlist.sourceforge.net/scowl-readme

5. Available at http://juilland.comp.lancs.ac.uk/dicer/

6. Available from the Oxford Text Archive: http://ota.ahds.ac.uk/

7. Variants here are words which VARD 2 deems to be variants, i.e. words which are not in its modern lexicon. It should be noted that words will be incorrectly marked as variants (particularly proper names) and some variants will be incorrectly marked as modern words (particularly *read-word errors*, such as 'bee' for 'be' and 'doe' for 'do').

# References

**Archer, D., McEnery, T., Rayson, P. and Hardie, A.** (2003). Developing an automated semantic analysis system for Early Modern English. In Archer, D, Rayson, P., Wilson, A. and McEnery, T. (eds), *Proceedings of the Corpus Linguistics 2003 conference*. UCREL technical paper number 16. UCREL, Lancaster University, pp. 22 - 31.

**Baron, A., Rayson, P. and Archer, D.** (forthcoming). Word frequency and key word statistics in historical corpus linguistics. To appear in *Anglistik. International Journal of English Studies, Special Issue: Corpus Linguistic – Methods and Approaches.*

**Culpeper, J. and Archer, D.** (forthcoming). The History of English Spelling. In Culpeper, J., Katamba, F., Kerswill, P., Wodak, R. and McEnery, T. (eds), *English Language and Linguistics*. Palgrave Macmillan, Basingstoke, UK.

**Görlach, M.** (1991). *Introduction to Early Modern English,* Cambridge University Press, Cambridge.

**Knapp, P. A.** (2000). *Time-Bound Words: Semantic and Social Economies from Chaucer's England to Shakespeare's*. Anthony Rowe Ltd, Chippenham, Wiltshire, UK.

**Kytö, M. and Voutilainen, A.** (1995). Applying the Constraint Grammar Parser of English to the Helsinki Corpus. *ICAME Journal* 19, pp. 23-48.

**Rayson, P., Archer, D., Baron, A., Culpeper, J. and Smith, N.** (2007). Tagging the Bard: Evaluating the accuracy of a modern POS tagger on Early Modern English corpora. In proceedings of *Corpus Linguistics 2007*, July 27-30, University of Birmingham, UK.

**Rayson, P., Archer, D., Baron, A. and Smith, N.** (2008). Travelling Through Time with Corpus Annotation Software. In Lewandowska-Tomaszczyk, B. (ed) *Corpus Linguistics, Computer Tools, and Applications – State of the Art. Palc 2007*. Peter Lang, Frankfurt am Main.

**Vallins, G. H., and Scragg, D. G.** (1965). *Spelling*. André Deutsch.