



Provided by the author(s) and University College Dublin Library in accordance with publisher policies. Please cite the published version when available.

Title	Review: Language and Chronology: Text Dating by Machine Learning (Toner and Han)
Authors(s)	Qiu, Fangzhe
Publication date	2020
Publisher	School of Celtic Studies, Dublin Institute for Advanced Studies
Link to online version	https://www.dias.ie/celt/celtica/celtica-contents/celtica-volume-32-2020/
Item record/more information	http://hdl.handle.net/10197/12243

Downloaded 2021-09-24T07:11:46Z

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



© Some rights reserved. For more information, please see the item record link above.

Language and Chronology: Text Dating by Machine Learning. Gregory Toner and Xiwu Han. Language and Computers 84. Brill, Leiden/Boston, 2019. xii + 183 pp. €88. ISBN 978-90-04-41003-9.

In 2015, two grants were awarded for projects using computational and statistical methods to date medieval Irish texts: one is provided by the European Research Council, which funds the project *Chronologicon Hibernicum* in Maynooth; the other is awarded by the Leverhulme Trust to fund Gregory Toner's project 'Dating of medieval texts through regressive analysis of the lexicon' in QUB. The present book is the outcome of the latter project, in which the two co-authors explore computational methods previously unknown to the field of medieval Irish studies and demonstrate the huge potential such methods embody for the discipline. When one compares the title of the project to that of the book, it is apparent that Toner and Han have pushed the methodological boundary much further by transcending from regression analysis in the initial project idea to advanced machine learning techniques in the outcome.

As the title of the series suggests, the book stands at the intersection of linguistics and computer science. Its main goal, as stated in the Introduction (pp. 2–3), is to achieve linguistic dating of medieval Irish texts by means of machine learning methods. Especially, the project adopts the supervised approach of machine learning, namely, human provides certain labels for the training data for the computer algorithm, which then learns a function that describes the patterns (linguistic features of a text) associated with each label (date). After cycles of training, evaluating and calibrating the algorithm, the computer programme can gradually improve its accuracy in performance (pp. 3–4), and thus 'learns' to do its task. Compared to human, the computer is capable of processing a large amount of data in a significantly shorter period of time, therefore it is possible to reveal association of patterns and labels that are hidden from human eyes.

The problem arising from using machine learning techniques is that, the actual programme that runs the algorithm is constantly shifting shape and often opaque even to humanities scholars with digital skills (pp. 7–8). What is more, such a 'black boxing' mechanism means that the actual patterns which the algorithm finally chooses to date a text are not visible (p. 8), so it is almost impossible to say how exactly are two texts similar or dissimilar linguistically. On the face of it, it seems as if the role of a linguist or a philologist has been reduced to one that passively accepts the date dictated by the machine, but this is actually not the case. In supervised learning, an algorithm performs as well as the training data that are fed to it, and only the humanities scholars can prepare clean, accurate data for the task. Expertise in linguistics, palaeography, history and other 'conventional' humanities domains is also indispensable for evaluating and interpreting the results (pp. 8–9). Therefore, close collaboration between humanities and computer science is necessary to engage in such an enterprise.

Accordingly, the book combines both the humanities and computational perspectives: on the one hand, there are discussions related to the nature of the medieval Irish texts, the challenge of dating them and the significance of the predicted dates, from the humanities scholars' perspective; on the other, detailed statistical and computational analyses of the data are also given, especially in chapter 2–4. Celtic linguists may be tempted to follow the advice on p. 9 to read only the introduction and conclusion of some of the chapters, but it will be beneficial to get oneself familiar with the basic methods (FTI, SVM, etc.) and parameters (n , t , δ , etc.) introduced in these chapters and mentioned throughout the book. Appendix B (pp. 162–164), which explains some fundamental concepts of machine learning, may be consulted before readers venture into the main chapters.

Chapter 1 'Dating texts: principles and methods' (pp. 11–40) offers a welcome reflection on conventional dating methods well known to scholars of medieval Irish texts, complementary

to Stifter (2013) which focuses on conventional linguistic dating parameters. These methods include dating by known authors, by internal evidence such as references to historical persons and incidents, by the date of the manuscript, by intertextual relationship (e.g. citation or borrowing between texts), by metrics and by linguistic features. All these methods, though useful to some extents and have yielded meaningful results, are beset with their particular problems. In my opinion, an ideal dating method should satisfy all of the three qualities: 1) granularity, that is, the date assigned should be as specific as possible, 2) applicability, that the method should be applicable to all targeted texts and 3) credibility, that the confidence of the assigned date should be high enough. However, for texts written in historical, under-resourced languages, such ideal method poses an impossible trinity, as is shown by the authors. For instance, while dating by known authors is granular enough, it is only applicable to a very small amount of medieval Irish texts, if we require a reasonable level of credibility (e.g. discrediting many poems that ‘Colum Cille recited’) (pp. 11–13). Dating by manuscript provides a credible *terminus ante quem*, and is widely applicable, but the granularity is so low that it often tells us next to nothing about the text, especially for early medieval Irish texts that are preserved in late medieval or early modern manuscripts (pp. 16–17).

While the limitations of other methods are mostly due to historical circumstances or the nature of the method itself, and therefore leave little chance for making further ground-breaking discovery, there is still room for developing the method of linguistic dating, as will be explained further down in the book. In this chapter the authors focus rather on identifying the possible pitfalls faced by the conventional linguistic dating method. The authors define the method as ‘the dating of a document or text based on the chronological stage of the language within it’ (p. 21). The authors’ description of this method on p. 23 and p. 38, however, is somewhat confusing. For instance, on p. 38 they observe that ‘[d]ating by linguistic means is normally approached in two ways. One is to observe the profile of linguistic features of known or accepted date...this directly provides an absolute date...The second approach is to sort texts chronologically through a comparison of their language overall.’ In fact, these are but one approach, as no such ‘absolute date’ for certain linguistic features exists independently, at least for medieval Ireland. The period in which a linguistic feature exists, such as the neuter gender, is defined by the texts in which such a feature is found or is used correctly, so after all it is the texts that are compared to each other.

It seems to me that the conventional dating method is best illustrated by the introductions to editions of early Irish texts. These usually compare the linguistic profile of the edited text to some landmark texts in different periods, explicitly or implicitly. The linguistic profile in comparison consists mostly of phonological and morphological features, but orthographical, syntactic and lexical features are also examined. The landmark texts are texts that can be more or less precisely dated, such as the Old Irish glosses in contemporary manuscripts, but previous scholarly opinions on the dates of other texts are also often quoted. Given a series of discrete landmark texts ($T_1, T_2 \dots T_n$), their respective periods ($t_1, t_2 \dots t_n$) and linguistic profiles ($P_1, P_2 \dots P_n$), the presumption is that linguistic changes are linear and proceed at a constant rate (cf. Kroch 1989), so that if the linguistic profile of a text P_i lies ‘between’ P_m and P_n , one can assume that T_i dates between T_m and T_n , but scholars sometimes propose a more precise t_i within half a century or less without giving specific reasons.

Obviously, there can be a lot of loopholes in this method, not the least because it has rarely been quantified and therefore relies heavily on learned guess. The authors point out other problems such as reliance on relative chronology (p. 23), linguistic strata and scribal revision (pp. 23–26), possible dialectal variation (pp. 26–27), register (pp. 27–29), deliberate archaism or obscurity (pp. 29–33). The discussions in these sections remind us that more typological studies are still needed on the interaction between texts, scribes and manuscripts, in order to use the correct dating methods for different circumstances. The ‘methodology’ section (pp.

36–39) reiterates the complicity arising from the typical transmission of early medieval Irish texts, and raises two principles in conventional linguistic dating (though not always followed in editions): 1) innovative features are more telling than conservative ones; 2) the proportions of features must be considered, in other words, quantitative measures must be taken. In general, the discussions offered in this chapter are valuable contributions to the topic, but not quite adequate. The chapter may benefit from recent scholarly literature on linguistic variation and change (e.g. Jensen and McGillivray 2017). However, since the aim of this book is to employ a different approach to the dating problem, this may not be necessary and I believe that a general reader will have a relatively clear idea of the status quo and challenges of dating early medieval Irish texts after reading this chapter.

Chapter 2 ‘Computational approaches to text dating’ (pp. 41–66) embarks on the more technical aspect of the project. The efforts described in this chapter exemplify nicely how to translate a research question from the humanities, viz. dating of texts, to a statistical and computational one. Text dating belongs to the larger domain of Natural Language Processing, a branch in vogue in computer science these days which underlies a number of day-to-day technologies around us. For a computer scientist, text dating can be translated as finding the most likely timestamp for a text based on the information provided by the text, which is regarded as a classification task that tries to categorise texts into classes defined as time intervals, or a regression task that ranks the texts according to temporal parameters (pp. 41–43). The selection of features to train the computer programme on is essential, as it informs the computer what and how features are relevant to the timestamp (p. 43). Most of the studies so far are based on lexical, semantic or topical features (p. 44), especially Named Entities, which are intuitively also how human usually recognises the period of a text (e.g. the mentioning of President Reagan or of Facebook). The following section (‘The Problem Stated’, pp. 44–47) translates the dating and evaluation task into mathematical language of functions and probabilities, which I refrain from repeating here. The authors then explain the three major types of solution for automatic text dating. The first is language modelling, which involves building a model of what the language looks like in each of the interval, mapping a text to the model, and training the programme to compute the function by which a new text can be dated (p. 48). This is basically a simulation of how linguists traditionally date a text by comparing selected features in the target text to those in the dated landmark texts, only augmented by the help of a machine. While this solution demands a huge amount of input from linguists, it offers ‘qualitative explanations about why a text is dated to a certain time, which is the most compelling aspect of the language modelling approach’ (p. 48) Such an approach is in fact employed by the ChronHib project in Maynooth. The second solution is the ranking approach, namely, to order a set of texts with respect to some measure (p. 49). The classification method regards timestamps as labels and, with suitable allowance of tolerance (labelled δ), the question becomes to which time interval of $2\delta+1$ does the text most likely belong (p. 50). In practice, these methods are sometimes combined to achieve the best performance, but the authors find that the multi-class classification could outperform the ranking approach (p. 51). The problem that remains for medieval Irish texts, of course, is whether and how the real temporal distribution underlying the training corpus can be modelled (p. 52). We will see this issue reflected in a later chapter. For the moment, the authors suggest five new solutions based on multi-class classification to improve the performance: a) Flexible time intervals (FTI), b) Sliding time interval (STI), c) Greedy grouping (GG), d) Temporal landmark selection (TLS) and e) FTI combined with TLS (pp. 52–64). These methods seek to find the best segmentation of time intervals.

One point to bear in mind for readers not familiar with Natural Language Processing is that in order to carry out these tasks, the texts have been treated as string vectors that consist of individual characters aligned in a one-way direction. These vectors are then analysed using

character *n*-grams that pick up 1, 2 or 3 characters at a time (p. 54). By doing so, the conventional control over feature selection, such as the division into phonological, syntactic or metrical features, is lost, and the algorithm cannot directly tell whether it is morphology, orthography or topic that triggers the classification. As the authors admit, the algorithm ‘is simply finding the most straightforward route to achieving maximum accuracy in the training texts’ (p. 133), and therefore all the ‘linguistic’ dating tasks referred to from this chapter on should be understood in the broadest sense of the word, that is, anything conveyed by language (in its written form). The algorithm, therefore, understands nothing of the content and is indiscriminate to any natural language. While this approach, again, creates a ‘black box’ that may displease philologists, it has the advantage of being strictly quantitative and unbiased towards a specific type of temporal feature, such as phonology in traditional editing and dating of early Irish texts.

Chapter 3 (‘Trials in English and medieval Irish texts’) (pp. 67–93) reports the testing results of the methods proposed in Chapter 2. The English texts tested here include historical newspaper excerpts between 1700 and 2010, and 180 days’ length of contemporary public posts from England. Each of the methods mentioned above has been used on these texts, with the reports attached. The medieval Irish texts tested are the Annals of Inisfallen (AI) 1092–1309 AD, the Annals of Ulster (AU) 1092–1378 AD, and the Annals of Loch Cé (ALC) 1014–1348 AD, harvested from the CELT digital editions. The reason of choosing the annals is that they provide a fine-grained periodisation down to a single calendar year, and a large part in AI data are contemporary records that supposedly reflect the real-time change in language (pp. 77–80). The datability test (pp. 90–92) also substantiates that AI is more datable, that is, the evolution process of the text is more regular, thus indirectly confirming the contemporaneity of AI 1092–1309. The result shows that the STI and FTI&TLS models perform the best in the trials, with the accuracy of over 60% for the annals.

While the predictive dating model built out of annalistic texts works well for the post-11th century annals, does it perform as well for other, longer Irish texts? In Chapter 4 (‘Dating long documents’, pp. 94–131), the authors set out to test the model again, firstly, on other sections of the annals, namely AU 500–1588 and ALC 1000–1652. The authors incrementally choose these annalistic texts that could be dated correctly by using the learned dating model, add them to the training corpus, and train new dating models on the extended corpus (pp. 95–96). By doing so with the FTI method, they achieve a dating performance from 36.43% when $\delta=3$ (tolerating a margin of error of 3 years), up to 84.76% when $\delta=50$ (p. 97). Within this extended corpus, the model is ineffective to predict the annals before 750 and only reaches reasonable accuracy (above 65%) after 1000 AD when $\delta=25$, that is, dating an annalistic entry to within 50 years. The reason for this may be due to the low number of syntactically continuous Irish texts before 750, and to the higher probability of contemporaneity in late medieval entries (pp. 96–98). The result is impressive, but, as the authors remind us on the high accuracy of annals 750–849, the model may be overfitted (p. 98).

The term ‘overfitting’ basically means that a model that fits a particular set of data very well may actually perform poorly on new data, because it depends too much on the noise of that dataset and contains too many parameters (Cohen and Jensen 1996). In order to test whether the model built upon annals is overfitted, the authors apply it to the dating of longer texts. They have selected 22 medieval Irish texts from the CELT corpus that cover the period c. 700 – c.1500. Each of the texts has a given date or range of date based on traditional methods (p. 99), and the authors run the algorithm on chunks of 20+ words in each text and take the most frequent estimated date as the final estimation. When the estimated dates are compared to the traditionally estimated dates, at a dating tolerance of 50 years, the accuracy is still at a poor 31.82% (pp. 100–101). In order to improve the performance, they select the top-3 most frequent predicts, and add in reinforcement and adaptive learning technique. The accuracy rate has then

been boosted to 40.63%. The improvement, however, is uneven, giving more accurate dates mostly for texts before 900 but creating larger margins of error for texts after that (pp. 102–109). The high accuracy achieved for 750–849 of the annals is not transferrable to longer, non-annalistic texts, and all such texts traditionally dated to before 850 are now given the estimated date of 914 by the machine (p. 111). The training model is obviously overfitted. Meanwhile, however, the model can successfully distinguish Early Modern Irish texts from Early Irish texts, and Middle Irish texts belonging to the 12th century are relatively well dated (p. 111–112). The different extents of robustness of the model in different periods reflect the quality of the training data: the period of 850–999 has fewer entries in average and lower accuracy in dating the annals, leading to fewer diagnostic features to enable correct dating, whereas the 12th century annals quite loyally reflect the contemporary language.

The same test is then extended to more texts that may have multiple strata due to their complex history of transmission. The accuracy rate reaches 42.86%, and the model again performs better within c. 900 – c. 1300 than the other periods (p. 113). The possible influence exerted by the manuscript date is also measured (pp. 118–119). The authors subsequently run the algorithm with a 10-year tolerance instead of 50, but only permitting the predicted dates within the 101-year period (± 50 years) determined by the previous test (pp. 119–121). This ‘focussing’ measure has generally shifted the proposed dates closer to those established traditionally, but, as shown by Figure 4.8 on p. 122, the dates given to Old Irish and Early Modern Irish texts still mismatch the traditional dates. After bias reduction, the focussed method is able to achieve a 50% accuracy rate, and generally succeeds in placing texts in the correct linguistic periods of Old Irish, Middle Irish and Early Modern Irish (pp. 123–125). For certain texts such as *Mescad Ulad*, where strata from different periods are juxtaposed in discrete chunks, the algorithm can even detect the strata by assigning different dates to the chunks. This will have interesting bearing on the study of textual history of some ‘composite’ texts, and I would be very eager to see if it works well in detecting layers in legal texts.

It has to be said that in this part, the numbered lists of texts on p. 101 and p. 113 have been merged and each text is given a new number, which could have easily confused the readers. For instance, ‘Text 14’ on p. 123 refers to *Saltair na Rann*, the date (872) assigned to which by the algorithm ‘is almost certainly too early’; the same tag, however, refers to *Cogadh Gaedhel re nGallaib* on p. 112, which is dated by the algorithm ‘squarely within the range suggested by scholars’. *Saltair na Rann*, on the other hand, is given the number A.1.30 in Appendix A ‘Conventional dating of texts used in this study’ (p. 156), which is arranged alphabetically. It would have been much clearer if the authors adhered to a single system of reference to texts.

The ‘Conclusion’ chapter (pp. 132–141) provides an overview of the contents of the book. In particular, the authors reflect on the important question of whether the model they develop is actually ‘linguistic’ (pp. 133–137). The algorithm probably relies on non-linguistic criteria such as topic, genre and Name Entities, as I have remarked above. The difference in the accuracy rate of dating the annals and the longer texts may, at least partly, be attributed to such non-linguistic factors. In the light of this, the authors call their model a ‘temporal model’, since the algorithm was fed with date labels and character vectors, and it is exactly the association between certain patterns in the character vectors and the date labels that it is predicting for us. In this way, the algorithm serves as a tool in computational chronometrics: it works very well in periodisation; it may help us detect strata in composite texts; and it has huge potential in dating medieval Irish texts with reasonable granularity, if we could use more data to boost and refine the capacity of the algorithm (pp. 137–140).

Appendix A details the conventionally assigned dates for the texts tested in the book, 36 in total plus three 11th- or 12th-century manuscripts (Lebor hUidre, Rawlinson B 502 and the Book of Leinster). This review is no place for examining all the arguments presented in this Appendix

or the dates accepted and used by the authors in their tests. I wish only to say that not infrequently a text has been assigned different and sometimes quite diverge dates by scholars, and the dates used by the authors should be regarded as tentative. Therefore, if the algorithm predicts a date that does not match the conventional date, the err may have rather occurred on the human side. For instance, although *Betha Adomnáin* is generally regarded as encoding political messages only relevant to the late 10th century and thus likely to be composed then (Herbert and Ó Riain 1988: 4–8), the algorithm assigns a date of 1120 (p. 120) which comes closer to the lifetime of Máel Ísa Ó Brolcháin (†1086), proposed by Mac Donncha (1976) as the author. The accuracy of the algorithm could indeed be higher than the authors think. The book concludes with another appendix on Machine Learning (pp. 162–164), a bibliography (pp. 165–180) and a subject index (pp. 181–183).

This book is the first to systematically apply machine learning method to the study of linguistic and literary history of medieval Ireland, and it opens up many possibilities and raise new questions. The algorithm, of course, is still at its initial stage and the accuracy cannot be compared to that achieved by commercial programmes run by Google or Facebook, but the methodology is sound and given more time and data, it could develop into a powerful tool that can significantly advance our understanding of the evolution of the Irish language and its textual tradition. The authors have done a beautiful job by combining humanities with computer science, in an accessible and informative way.

References:

- Cohen, Paul and Jensen, David 2000: ‘Overfitting Explained’. *Preliminary Papers of the Sixth International Workshop on Artificial Intelligence and Statistics*.
- Herbert, Máire and Ó Riain, Pádraig (eds) 1988: *Betha Adomnáin: the Irish Life of Adomnán* (London).
- Jenset, Gard B., and McGillivray, Barbara 2017: *Quantitative Historical Linguistics: A Corpus Framework* (Oxford).
- Kroch, Anthony 1989: ‘Reflexes of grammar in patterns of language change’, *Language Variation and Change* 1 (3), 199–244.
- Mac Donncha, Frederic 1976: ‘Medieval Irish Homilies’, in Mairtín McNamara (ed.), *Biblical Studies: The Medieval Irish Contribution* (Dublin), 59–71.
- Stifter, David 2013: ‘Towards the Linguistic Dating of Early Irish Law Texts’, In Anders Ahlqvist and Pamela O’Neill (eds) *Medieval Irish Law: Texts and Contexts* (Sydney), 163–208.

Fangzhe Qiu
University College Dublin