



Universiteit  
Leiden  
The Netherlands

## Contributions to the computational processing of diachronic linguistic corpora

Cunha, E. Landulfo Teixeira Paradela

### Citation

Cunha, E. L. T. P. (2020, March 19). *Contributions to the computational processing of diachronic linguistic corpora*. LOT dissertation series. LOT, Amsterdam. Retrieved from <https://hdl.handle.net/1887/133504>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/133504>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/133504> holds various files of this Leiden University dissertation.

**Author:** Cunha, E. Landulfo Teixeira Paradela

**Title:** Contributions to the computational processing of diachronic linguistic corpora

**Issue date:** 2020-03-19

## CHAPTER 3

---

### Establishment and obsolescence of linguistic items in a diachronic corpus<sup>1</sup>

---

#### 3.1 Introduction

Diachronic and historical corpora are useful tools to study linguistic phenomena that unfold over time, including processes of variation and change. Previous work has employed these kinds of corpora to analyze language change in progress (Hundt and Mair, 1999), to infer cases of variation and change (Bauer, 2002), and to investigate language change using word vector embeddings (Hamilton et al., 2016), to mention a few.

When using diachronic corpora for investigating language variation and change, one of the relevant tasks for researchers is the iden-

---

<sup>1</sup> This chapter reproduces with minor changes the article “An algorithm to identify periods of establishment and obsolescence of linguistic items in a diachronic corpus” (Cunha and Wichmann, in press), accepted for publication in *Corpora*. See Appendix C for more information.

tification of specific time periods in which certain linguistic items arise and, conversely, vanish. It is particularly valuable to detect when items (i) are first attested, (ii) become established in the corpus, (iii) become obsolete and (iv) are attested last. Although the detection of the earliest and the latest attestation dates of items in a diachronic corpus is trivial, the same cannot be said about their establishment and obsolescence, because there are no clear and commonly accepted criteria for pinpointing when an item is getting established and when it can be regarded as obsolescent (cf. Tichý, 2018).

The aim of this chapter is threefold: first, to formulate a set of criteria to define binary notions of establishment and obsolescence of items in a diachronic corpus; second, to present an algorithm to assist in the identification of specific time periods of establishment and obsolescence of linguistic items in diachronic corpora according to the previously mentioned criteria; and, finally, to use this algorithm to make a series of general considerations based on real data for the purpose of demonstrating the utility of the methods presented here and for making some observations on two centuries of the dynamics of the American English lexicon that are interesting in their own right. We will observe, among other findings, that the proportion of words established in a given decade is similar across decades and, by studying the words stemming from different decades that are most frequent today, we will get an impression of how the lexical heritage of contemporary American English bears the imprints of salient aspects of life as it was experienced during specific, previous decades.

The algorithm proposed here is simple and generalizable. It can be applied to any corpus that is divided into time frames, regardless of language or historical period, since it only takes as input information on the frequency of the analyzed items in each time frame. Likewise, the nature of the items under analysis is, in principle,

irrelevant to the applicability of our algorithm, so it can also be implemented to examine aspects of language not considered in our case studies, such as phonology or morphosyntax. Moreover, the algorithm, or some derived version, should be generally applicable to the investigation of time series of sociological, anthropological or historical data.

### 3.1.1 Related work

Previous quantitative investigations on language dynamics have dealt with the notions of birth and death of linguistic items, which are related to the concepts of establishment and obsolescence contemplated here. Petersen et al. (2012), for instance, analyze more than 200 years of data from three different languages with the goal of shedding light on the aggregate dynamics of word evolution in written texts. They investigate variations in the use of words during their lifespans and, among other results, identify a tendency for a peak in word use growth rate to occur around 30-50 years after a word's first attestation in their corpus. Furthermore, the authors find evidence that the dynamics of word evolution might be influenced by historical events, such as wars. This last observation is also made by Bochkarev et al. (2014), who additionally find a relationship between the frequency of a word and its stability in the lexicon of a language, confirming previous results from Pagel et al. (2007). Moreover, Perc (2012) analyzes the evolution of high-frequency English words and phrases, discovering that their lifespan is not uniform across the centuries, and Michel et al. (2011) investigate some patterns in the evolution of English lexicon and grammar. In addition, the work performed by Kerremans et al. (2012) in the scope of the NeoCrawler project<sup>2</sup> presents a Web crawler that iden-

---

<sup>2</sup>Available at <http://www.neocrawler.anglistik.uni-muenchen.de/crawler/html/>.

tifies and retrieves neologisms from the Internet, making it possible to analyze “how words spread in the early stages of their life cycles and how they begin to establish themselves in lexical and semantic networks” (p. 59).

Certainly connected to the concepts of first attestation, establishment, obsolescence and last attestation of an item in a corpus are the studies that use diachronic corpora to investigate language variation and change. Biber and Gray (2011), for instance, analyze the influence of written language on grammatical change, and suggest that new grammatical uses and functions emerge not only in spoken interaction, but also in written registers. Topics such as the variation of the English genitive (Hinrichs and Szmrecsanyi, 2007), the variation of complex prepositions in Brazilian Portuguese (Shepherd, 2014) and the change in the grammar of English verbs (Hilpert and Mair, 2015), to illustrate, have been considered in previous investigations that made use of diachronic corpora. The use of corpora with the aim of investigating creativity in literary and ordinary language, including novel word formation, is scrutinized by Vo and Carter (2010), while Moon (2010), in tackling the question of what corpora can reveal about lexicon, mentions that these tools might contribute to the analysis of the establishment and the institutionalization of new derivations and compounds in a language.

The notion of establishment of linguistic items in a diachronic written corpus from a particular language is not to be confused with the concept of entrenchment of structures in the memory of speakers (Langacker, 1987), which is central in the field of cognitive linguistics. Nevertheless, Schmid (2007) considers that this notion of entrenchment “also applies to language as such and whole speech communities, because the frequency of occurrence of concepts or constructions in a speech community has an effect on the frequency with which its members are exposed to them” (p. 119). As a consequence, it should be possible to talk of a degree of entrenchment

of a linguistic item not only in the memory of individual speakers, but also in a specific language. Indeed, Croft (2000) uses the notion of entrenchment in his proposal of an evolutionary model of language change, advocating for a strong relationship between the perpetuation of a given linguistic structure in the language and the degrees of entrenchment of this particular structure in the grammar of speakers. However, in discussing the relationship between frequency in natural language use and the entrenchment of complex linguistic strings in the minds of language users, Blumenthal-Dramé (2012) argues for a weak version of the so-called *corpus-to-cognition principle* – since, according to her, only “certain corpus-extracted variables may, to some extent, be used as a yardstick for entrenchment in the brain of an average language user” (p. 205). In this study, it is not our goal to contemplate entrenchment in the memory of speakers, nor to elaborate on the relationship between the frequency of linguistic items in a corpus and their entrenchment in the minds of individuals. For this reason, we opted for the use of the term *establishment* and, by not using the loaded term *entrenchment*, we hope to avoid any kind of misinterpretation of the goals of our proposed method.

Regarding the opposite phenomenon – that is, the loss of linguistic items –, Tichý (2018) presents one of the few studies on lexical obsolescence and mortality in English. Using a fine-grained methodology based on the difference between frequency levels in distinct periods of time, the author proposes a method for extracting from large corpora forms that were once common but later became obsolete. Our methodology differs from his in that Tichý is mostly interested in words that were once very common in the language, while our proposed methodology is more flexible in this regard. Also, Tichý’s proposal, being more fine-grained, is more computationally demanding, whereas our methodology is simpler and more straightforward. We consider the approaches complementary and

imagine that they may even be used together in some specific situation.

Finally, the work of Hilpert and Gries (2009) provides several resources for the assessment of frequency changes in multistage diachronic corpora. The authors present suggestions for the analysis of this kind of data, displaying examples and use cases of great value for historical linguists. In particular, we mention the introduction of the *iterative sequential interval estimation* (ISIE), a method that provides a range of expected frequencies for an item in each time period of the corpus. When the frequency “happens to go beyond the expected values, we have detected a change that merits further attention” (Hilpert and Gries, 2009, p. 393).

### 3.2 Defining establishment and obsolescence as binary notions for diachronic corpus linguistics

Dictionaries and glossaries of neologisms (e.g. Ayto, 1989, 1990, 1999; Tulloch, 1991; Algeo and Algeo, 1993; Knowles and Elliott, 1997, to mention works on English) attempt to record recent additions to the language, but their editors are usually aware that what they characterize as *new words* might not be new at all. In fact, Tulloch (1991) mentions the potential gap between the point in time in which a word enters the language and the moment when the general public becomes aware of it – which is the occasion when the neologism might be included in the most prestigious dictionaries and can be considered “established in the language” (Ayto, 1999, p. iii). Still, most of the past studies mentioned in Section 3.1.1 that analyze time periods in which linguistic items arose and vanished associate birth and death with, respectively, first and last attestations in a corpus. In this study, we argue and show evidence that



the first appearance of an item in a corpus may occur considerably earlier than its establishment in the corpus itself and, conversely, that an item might still appear in the data long after it became obsolete (see Section 3.4). This fact suggests that it may often be convenient to discriminate between first attestation and establishment as well as between last attestation and obsolescence, so as to obtain a more accurate description of the lifespan of a linguistic item.

As pointed out by Widdowson (2000), it is important to emphasize that a corpus is different from a language and, consequently, that the establishment or the obsolescence of an item in a corpus does not necessarily imply its establishment/obsolescence in a language. At most, it might be claimed that a corpus represents part of a language and that a relationship between these two entities exists.

We are interested in defining binary (rather than continuous<sup>3</sup>) notions of establishment/obsolescence in order to indicate whether a linguistic item may have arisen in or vanished from a diachronic corpus during the period covered by it. This is particularly useful for researchers interested in extracting lists of candidate items for further research (see Section 3.4.3). We stipulate that, in a particular corpus which includes diachronic information, each linguistic item (be it a word, a morpheme, a syntactic structure or other) may usefully be classified as being in one – and only one – of the following possible states in a given period: (a) established; (b) obsolete; (c) permanent; (d) short-lived; (e) random. These states refer to diachronic patterns of appearance of the item through the corpus. The state *established* concerns items that, although not frequent (above a given threshold) in the beginning of the period, rise in

---

<sup>3</sup> In other words, our aim is to provide sets of (candidate) established/obsolete items rather than some sort of “degree of establishment/obsolescence” per item.

frequency at some point and remain frequent until the end of the period covered the corpus. In other words, established items were not part of the language represented by the corpus, but at some point during its time span they flourished and remained frequent afterwards. The state *obsolete*, conversely, refers to items that are frequent (above a given threshold) in the beginning of the period covered by the corpus, but which at some point decrease in frequency. They are, therefore, items no longer in general use by the end of the corpus, although they may linger on as old-fashioned forms or archaisms making occasional appearances. The state *permanent* describes items that are frequent enough through the whole period covered by the corpus. The state *short-lived* regards items that flared up for some time and then, still during the period covered by the corpus, decreased in frequency again. Finally, the state *random* is reserved for items that do not show any of the aforementioned patterns. In the next section, we further develop this categorization by presenting our proposed methodology for classifying items into the above-mentioned classes.

## 3.3 The algorithm

### 3.3.1 Requirements

In order to be accessed by our proposed algorithm, a corpus must be divided into time frames. These time frames might delineate any desired period of time, depending on the nature of the data and on the research goals. Each one of these time frames may represent, for instance, a period of several years, or one decade, or one year, or even one day – the latter in the case of research using data from online social media platforms, for example. For methodological reasons, it is to be preferred that time frames are uniform (both in corpus size and duration, whenever possible) across the whole

corpus, but this is not a strict requirement and alternative methods (such as the one proposed by Gries and Hilpert (2008)) could be used to divide the corpus in time stages. Also, our method relies on the use of topically coherent corpora, so as to avoid that changes in sampling across time lead to change in the frequency of linguistic items.

In our method, when the frequency of a given item in a certain time frame is above a definite threshold, it is represented by the digit 1; when this frequency is below this threshold, by 0. For example, in a corpus divided into six time frames, the *diachronic sequence* of an item whose frequency exceeds the threshold only in the last time frame is denoted by 000001, while the sequence of an item whose frequency exceeds the threshold in all but the second and third time frames is denoted by 100111.

We leave the definition of the boundary between assigning a 0 or a 1 in the diachronic sequence as a choice for the researcher who will use our algorithm, since this depends on additional methodological choices and assumptions. We strongly discourage, however, the use of absolute frequencies as thresholds (as they are dependent on the size of the corpus in each time frame) and, conversely, encourage the use of relative frequencies. For example, a 1 might be attributed to a given item in a particular time frame in case its frequency exceeds  $n\%$  of the total size of the corpus in that time frame; otherwise, a 0 will be attributed. A simple and useful case is when this boundary is set on a really low relative frequency (e.g. 0,00000001% of the corpus size). In this case, the mere presence of the item in the time frame is enough to assign a 1 to it. This simple situation is convenient, practical and might still give interesting results, such as the ones we display on Section 3.4.

In Section 3.3.2, we introduce the rules regulating a first algorithm aimed at the categorization of linguistic items into one of the previously mentioned states – *established*, *obsolete*, *permanent*,

*short-lived* or *random*. We begin by stating naive rules that are ultimately not satisfactory for our intentions. In Section 3.3.3, however, an improved version of these rules, more effective for the purposes of the goals declared here, is presented.

### 3.3.2 Rules for a naive algorithm

A first (and naive) version of an algorithm aiming to solve the task of categorizing a linguistic item into one of the aforementioned states may be based on the following rules:

- Established items: those that are not frequent enough in the corpus before a certain time frame, but from a given point start to exceed the frequency threshold in all of the following time frames, without exception. Example of a diachronic sequence in a corpus containing six time frames: 000111.
- Obsolete items: those that are frequent in the first time frame(s), but from a given point onwards are not frequent enough in any of the following time frames, without exception. Example: 111000.
- Permanent items: those that are frequent in all time frames, without exception. Only possible diachronic sequence: 111111.
- Short-lived items: those that are not frequent enough in the extremes of the period covered by the corpus, but that are consistently frequent during an intermediate period. Example: 00111100.
- Random items: those that do not fit into any of the previous cases. Example: 100101.

It is clear that these rules only work for what we might call “perfect” patterns, in which linguistic items “appear” or “disappear” at a certain point and keep this status until the end of the period covered by the corpus, without fluctuations. According to this method,

an item which, in a corpus divided into ten time frames, exhibits the pattern 0001011111 is considered an example of a random pattern, even though it is obvious for us that it clearly illustrates an item established sometime around the middle of the period covered by the corpus. To solve this issue, an improved version of these rules, allowing for some deviations from perfect patterns, is presented in the next section. Without the allowance of these deviations, the low frequency of an item in a specific time frame would be too severely punished, being enough to disregard the item as an innovation; conversely, the presence of an item in a specific time frame could be enough to disregard it as an obsolete item.

### 3.3.3 Proposed algorithm

Here, we propose an algorithm that enhances the previous approach by allowing for small deviations from perfect patterns, thus making it possible to include more (and more accurate) data into the lists of established and obsolete items of a corpus. The core idea is (i) to compare the observed (real) diachronic sequences of each item in the corpus with perfect patterns for establishment and obsolescence, and then (ii) to select a specific time frame as representing the time of establishment or obsolescence, using the criterion that it should be the time frame that produces the smallest amount of deviation from these perfect patterns.

Consider the following fictitious example. In a corpus divided into ten time frames, the linguistic item *A* exhibits the diachronic sequence 0001110111 – according to which *A* is not frequent enough in the initial periods of the corpus, but after time frame four it is consistently frequent, with the only exception being time frame seven. Our algorithm inspects each position in between two adjacent time frames, starting from position one (which lies in between the first and the second time frames, as in 0\_001110111). The perfect

pattern indicating the establishment of an item in this position is 0\_111111111 (i.e., the item is not present before the position and is consistently present after it), while the perfect pattern indicating its obsolescence at this point is 1\_000000000. Here, the algorithm investigates the observed sequence for item *A* and counts deviations from the two perfect patterns. By *deviations* we mean differences in particular points of the diachronic sequences: for instance, if, in a given place, a 0 is found in the observed sequence when a 1 is expected according to the perfect pattern, then we detect a deviation<sup>4</sup>.

Let us return to the example of the sequence 0001110111. At the first position, when the assumption is that the item gets established after that point in time, the algorithm finds three deviations from the perfect pattern (the three 0s in time frames two, three and seven); when the assumption is that the item becomes obsolete, there will be seven deviations from the perfect pattern (the 0 in time frame one and the six 1s in time frames four, five, six, eight, nine and ten), as illustrated below, where arrows indicate deviations:

Observed sequence	0_001110111		Observed sequence	0_001110111
	↓ ↓ ↓			↓ ↓ ↓ ↓ ↓ ↓ ↓
Perfect pattern (establishment)	0_111111111		Perfect pattern (obsolescence)	1_000000000

After these results have been obtained for the first segmentation, the algorithm moves to the next position (00\_01110111). Here, two deviations from the perfect pattern of establishment (the two 0s in time frames three and seven) and eight deviations from the perfect pattern of obsolescence (the two 0s in time frames one and two, and the six 1s in time frames four, five, six, eight, nine and ten) are

---

<sup>4</sup> These deviations might be counted, for example, by employing an edit distance algorithm, such as the Levenshtein distance algorithm, that returns the minimum number of single-character edits required to change one sequence into the other.

found. In the third position (000\_1110111), only one deviation from the perfect pattern of establishment is found (the 0 in time frame seven), while nine deviations from the perfect pattern of obsolescence are detected (the three 0s in time frames one, two and three, and the six 1s in time frames four, five, six, eight, nine and ten). In the next step, two deviations from the perfect pattern of establishment (the 1 in time frame four and the 0 in time frame seven) and eight deviations from the perfect pattern of obsolescence (the three 0s in time frames one, two and three, and the five 1s in time frames five, six, eight, nine and ten) are identified in the fourth position (000\_1110111). The process continues until all positions<sup>5</sup> are analyzed, after which the position producing the smallest number of deviations can be found. This position will represent a possible moment of establishment or obsolescence. In the case of item *A*, Table 3.1 shows that the smallest number of deviations is found under the assumption of establishment (rather than obsolescence) and is observed in position three, indicating that this linguistic item might have been established in the corpus immediately after this point – that is, within time frame four.

Let us go on to consider, in the same corpus, a linguistic item *B* exhibiting the diachronic sequence 1111100100. After the inspection of the nine positions in between each two adjacent time frames, the proposed algorithm outputs that the smallest number of deviations from a perfect pattern is found in position five, but now the assumption is that of obsolescence. In this case, the decision implies that the item has become obsolete in the time frame following that position, which corresponds to time frame six, as displayed again in Table 3.1.

---

<sup>5</sup> Note that the number of positions to be analyzed equals  $tf - 1$ , where  $tf$  represents the number of time frames in which the inspected corpus is partitioned.

Table 3.1: Number of deviations in each position according to the proposed algorithm for two fictitious examples  $A$  and  $B$ . In the first example, the smallest number of deviations is observed in position three under the assumption of establishment, indicating that  $A$  might have gotten established immediately after that position (in time frame four). In the second example, the smallest number of deviations is observed in position five under the assumption of obsolescence, suggesting that  $B$  may have become obsolete immediately after that position (in time frame six).

Item $A$			Item $B$		
Position in the diachronic sequence	Deviations (establishment)	Deviations (obsolescence)	Position in the diachronic sequence	Deviations (establishment)	Deviations (obsolescence)
1 (0_001110111)	3	7	1 (1_111100100)	5	5
2 (00_01110111)	2	8	2 (11_11100100)	6	4
<b>3</b> (000_1110111)	<b>1</b>	9	3 (111_1100100)	7	3
4 (0001_110111)	2	8	4 (1111_100100)	8	2
5 (00011_10111)	3	7	<b>5</b> (11111_00100)	9	<b>1</b>
6 (000111_0111)	4	6	6 (111110_0100)	8	2
7 (0001110_111)	3	7	7 (1111100_100)	7	3
8 (00011101_11)	4	6	8 (11111001_00)	8	2
9 (000111011_1)	5	5	9 (111110010_0)	7	3



It is worth noting that, for each position, the number of deviations from the perfect pattern indicating establishment plus the number of deviations from the perfect pattern indicating obsolescence equals the amount of time frames in the corpus. This is obviously expected, since each 0 or 1 in the observed sequence is always a deviation from a perfect pattern (either regarding establishment or obsolescence), but never a deviation from both.

The proposed algorithm will always output a smallest number of deviations from the perfect patterns, but this value might be considered excessive in some cases. For this reason, a cut-off point of the number of acceptable deviations from establishment and obsolescence should also be defined, and cases that exceed this threshold should be assigned to the pool of cases of random distributions. This cut-off point must be set by the researcher according to some sensible considerations that will vary according to the type of corpus in question: if it is a lexical corpus of child language acquisition with day-to-day recordings, for example, there might be many deviations since a single child is not expected to exercise its full vocabulary every day; if it is a large historical corpus of texts with yearly time frames, the cut-off point could be set to fewer deviations<sup>6</sup>. Here we must refrain from generalization about such thresholds, but we give an example of how to derive one from the behavior of a specific corpus in Section 3.4. What we mainly want to stress is that the use of an algorithm such as the one described here has the advantage that there has to be such an explicit threshold. Even if it is defined in somewhat *ad hoc* ways in individual cases, it will force researchers to be specific about their choice, enhancing transparency and replicability of a given study.

Finally, in case of ties such that the smallest number of devia-

---

<sup>6</sup> It is trivial to observe that the naive rules presented in Section 3.3.2 correspond to the algorithm proposed here when this cut-off point equals to zero.

tions occurs at more than one position, we advocate for choosing the position that includes more time frames with the item being present so as to maximize the amount of positive attestations after minimizing the amount of deviations. An example would be as follows: in a diachronic sequence such as 1111101000, where the smallest number of deviations from a perfect obsolescence pattern (which is one) is achieved both in positions five and seven, we favor choosing the latter (corresponding to the eighth time frame) as the moment of obsolescence; conversely, in a sequence such as 0001011111, we favor choosing the fourth time frame (rather than the sixth) as the moment of establishment.

In conclusion to the present section, we present a summary of the steps made by the algorithm.

**Summary of the algorithm:**

1. Go to the first position in between two adjacent time frames.
2. Calculate the number of deviations from the perfect patterns of both establishment and obsolescence.
3. If there are unexplored positions in between two adjacent time frames, go to the next position and repeat step 2; otherwise, go to the next step.
4. Compare the value found for the smallest number of deviations  $S$  with the maximum threshold for deviations allowed  $T$ ;
  - I. If  $S > T$ , the item is considered neither established nor obsolete.
  - II. If  $S \leq T$ :
    - i. resolve potential ties by choosing the position that includes most time frames with the item being present;
    - ii. consider the time frame immediately after the corresponding position as the time frame of establishment or obsolescence.

The previously described algorithm is able to identify items classified as *established* or *obsolete* according to our defined criteria, but not items evaluated as *short-lived* – which are classified as *random* by it. In Section 3.4.5, we provide a case study in which we suggest a way of adapting our method for this specific situation.

In the next section, we apply our algorithm to a real corpus, supplying five case studies to illustrate its usage and some of its potential for producing interesting observations.

### 3.4 Case studies

In order to demonstrate the applicability of the algorithm proposed in Section 3.3.3, we applied it to the Corpus of Historical American English (COHA). This corpus contains more than 100,000 texts from different sources (fiction and non-fiction books, magazines, and newspapers) published in the United States of America from 1810 to 2009 (Davies, 2012), and can be explored online and downloaded from its webpage<sup>7</sup>. In this work, we use the case-insensitive list of unique words<sup>8</sup> (types), annotated with part of speech (PoS) tags. This list contains the frequency of each pair (word + PoS tag) in each of the twenty decades spanned by the data. In this way, it is often possible to differentiate between homonyms (e.g. *light*, that can be tagged as adjective, noun, verb and others). We also removed all words classified with the tags for “formula”, “proper noun” (neutral for number, singular and plural), “letter of the alphabet” (singular and plural), “foreign word” (such as *arbre*, *bueno* and *deum*) and “unclassified word” (which includes ideophones like *bang-bang*, unrecognizable words such as *carige*, exclamations like

<sup>7</sup> <https://corpus.byu.edu/coha/>

<sup>8</sup> Here, we define a *word* simply as a string of characters uninterrupted by a space. It deserves mentioning that the downloadable COHA frequency data excludes words that occur less than three times in total in the corpus.

*gotcha* and recognizable words whose context is apparently unexpected). In total, we analyze 381,698 pairs of word + PoS tag in this corpus. As mentioned in Section 3.3.1, in these case studies we set the boundary between a 0 and a 1 in a really low relative frequency (0,00000001% of the corpus size) – so, the mere presence of a word in a time frame is enough to assign a 1 to it. By using this straightforward criterion, our goal is to show that even a method based on the simple presence/absence of items in specific time frame is able to rapidly bring interesting and useful results.

Having selected the corpus to work with, we need to decide on the value of  $T$ , i.e., the maximum threshold for how many deviations from the perfect patterns we can accept so to advocate for the establishment or the obsolescence of the analyzed items. Although, as mentioned in Section 3.3.3, the decision must to some degree be *ad hoc*, it should at least be backed up by an explicit criterion. Our approach here is to look at the statistics of establishment of words using the perfect pattern (no deviations) as a baseline: if the number of words that get established in different decades allowing for  $d$  amount of deviations is consistently proportional to the number of words that get established under the zero-deviation criterion, then the given value of  $d$  is acceptable. But how should “consistently proportional” be defined? Here, we look at the time series for the proportion of words that became established in each decade out of all words in the decade using different values of  $d > 0$ , and correlate these numbers with the corresponding numbers for the zero-deviation curve. If the  $p$ -value of a Pearson correlation is below 0.05 for a given value of  $d$ , then that amount of deviation is taken to be acceptable. In our case, it is only for  $d = 1$  that we find an acceptable correlation:  $p = 0.0047$ ,  $\rho = 0.605$ ; for  $d = 2$  we already get  $p = 0.0569$  and the correlation goes down to  $\rho = 0.4323$ . Results continue to get worse as more deviations are allowed for. Thus, it is clear that too much noise would be admitted into any

statistics on the establishment of new words (and presumably on their obsolescence as well) if more than one deviation is considered acceptable in this case. For one deviation, the observations will also contain some noise, but more (and still reliable) data will be included<sup>9</sup>.

As an illustration of how a few words are evaluated by our algorithm in COHA, Table 3.2 displays the outcomes of attempts to detect established/obsolete words using respectively the naive (zero-deviation criterion) approach and our proposed method implementing the one-deviation criterion. The words selected for illustration are all singular common nouns present in COHA. Words are marked with (a) when they represent cases in which their first/last attestation matches the outcomes of both algorithms; with (b) when the naive rules cannot determine their date of establishment/obsolescence and our algorithm finds that the first or last occurrences are, respectively, also the decades of establishment or obsolescence; with (c) when the naive rules again cannot tell their date of establishment/obsolescence and our algorithm now finds that the first or last occurrences are, respectively, *not* the decades of establishment and obsolescence; with (d) when the decade of establishment/obsolescence is considered random by both methods. The (b) and (c) cases are particularly relevant since they illustrate data that would be lost from the purview of a study of lexical establishment or obsolescence if no deviations were admitted.

---

<sup>9</sup> We stress that the decision on the value of this maximum threshold of deviations allowed must to some degree be *ad hoc*: since there are no “right” and “wrong” sets of established/obsolete items, this threshold depends on whether the researcher desires to obtain more comprehensive lists or more restricted ones – for the former, a higher threshold could be stipulated; for the latter, a lower value should be set. The point about using correlations and the  $p$ -value is that the distribution with one deviation from the perfect pattern is significantly similar to a distribution without any deviation, so the deviation can arguably be ignored.

Table 3.2: Outcomes of attempts to detect established/obsolete words using a first/last attestation approach, an algorithm following naive rules and the proposed algorithm (with a one-deviation criterion) in a selection of words present in the Corpus of Historical American English (COHA). Each time frame represents a decade, ranging from 1810s to 2000s. The examples chosen are all words tagged as “singular common noun”.

Word	Observed diachronic sequence	First attestation	Outcome according to	
			naive rules	proposed algorithm
(a) <i>victrola</i>	00000000001111111111	1910s	established (1910s)	established (1910s)
(b) <i>snatcher</i>	00000000110111111111	1890s	random	established (1890s)
(c) <i>bulldozer</i>	00000010000011111111	1880s	random	established (1940s)
(d) <i>wife-murderer</i>	00001011101100010010	1850s	random	random
Last attestation				
(a) <i>secrecy</i>	11111111111100000000	1920s	obsolete (1930s)	obsolete (1930s)
(b) <i>gratulation</i>	11111111111110100000	1960s	random	obsolete (1970s)
(c) <i>destruction</i>	11100001000000000000	1880s	random	obsolete (1840s)
(d) <i>unfeelingness</i>	00110000110110100100	1980s	random	random

### 3.4.1 Case 1: Statistics on established and obsolete words

Figure 3.1 shows the percentage of words that became established (left figure) and obsolete (right figure) per decade in COHA according to our algorithm and using the one-deviation criterion. In the left figure, the U-shaped nature of the curve concerning the establishment of words considering the whole corpus is easily explained by two factors that must always be acknowledged by the researcher: first, the proportion of words that had not appeared previously in the corpus is necessarily higher in the first time frames than in the next ones, as a consequence of the phenomenon known as Herdan’s or Heaps’ law (Herdan, 1964; Heaps, 1978), according to which vocabulary size grows slowly compared to the size of the document/corpus; second, the proportion of words arisen in a certain decade that are consistently present in the following ones (i.e., the words considered established conforming to our criteria) is necessarily higher in the last time frames than in the previous ones, because most of these recently established words did not have time to become obsolete yet. To demonstrate these two effects more precisely, we include two additional curves in the graph, corresponding to the percentage of words that became established in a given decade considering only certain time windows (six- and eight-decade windows). In other words, we reduce the corpus to sliding windows of six and eight decades in order to decrease the “advantage” that early and late decades hold compared to middle decades. For these two additional curves, however, we are employing a zero-deviation criterion, since one deviation in a universe of only a few decades might be considered disproportionate. These additional curves do not display such a clear U-shaped nature, even though the one regarding six-decade windows still slightly reflects this pattern especially in its left tail. Also, both exhibit the same shape, suggesting that the

proportion of established words among all words in a given decade is similar across time and that the use of different windows in this case might be no more than a question of how much data one wishes to consider: around 3% for six-decade windows vs. around 2% for eight-decade windows.

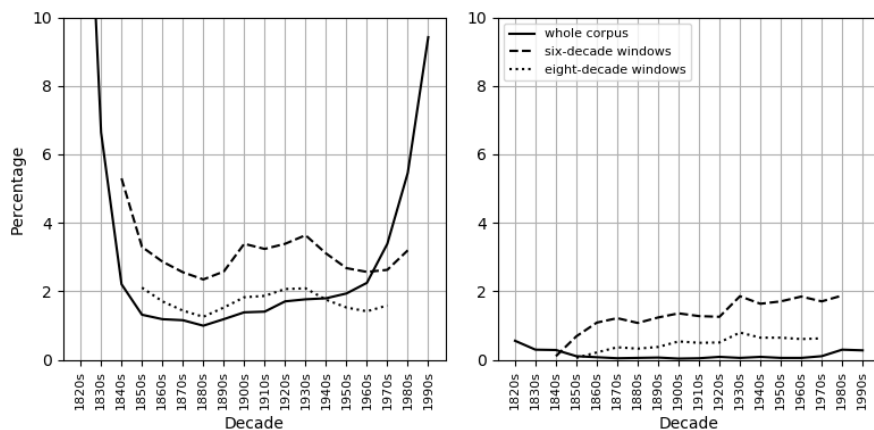


Figure 3.1: Percentage of words that became established (left) and obsolete (right) per decade using a one-deviation criterion and applying six- and eight-decade windows combined with a zero-deviation criterion. Curves comprise different time spans according to the sizes of the sliding windows.

Regarding the right figure, we observe that the proportion of words that became obsolete among all words in a particular decade is also more or less constant, with lower percentages than the ones referring to established words. Additional research must be carried out to more precisely understand the meaning of these results and their implications for language dynamics.



### 3.4.2 Case 2: Characteristics of established and obsolete words

Investigating certain characteristics of words considered established or obsolete according to our proposed algorithm is also a possible line of study. Figure 3.2, for example, shows the average length (in number of characters) of words that became established and obsolete in given decades. Here, we observe an irregular shape of the curve concerning words that became obsolete, but a consistently positive slope in the curve regarding established words, indicating a persistent increase in the average length of words established in the corpus across time – that goes from around eight and a half characters in the mid-19th century to almost ten characters in the second half of the 20th century. Since COHA is balanced by genre across time (Davies, 2012), this finding should in principle not be attributed to artifacts of the corpus (such as a potential increase in the proportion of scientific literature, for example). Additional investigation should be conducted to better understand this phenomenon, presumably employing other data and associating this results with the abundant previous work on word length (Grzybek, 2007).

Different analyses can be carried out also considering the PoS tags of the words in the corpus. Figure 3.3 depicts the percentage of parts of speech (grouped as “adjective”, “adverb”, “noun”, “verb” and “other”) among words that became established (left figure) and obsolete (right figure) in each decade. Among the established words, we visually notice a descending trend in the proportion of verbs and an ascending trend in the proportion of adjectives across the decades. The other curves are not consistently rising or falling – although, if we consider only the time period starting in the 1960s, we do observe a tendency for the proportion of nouns among the established words to increase. Regarding the words that became

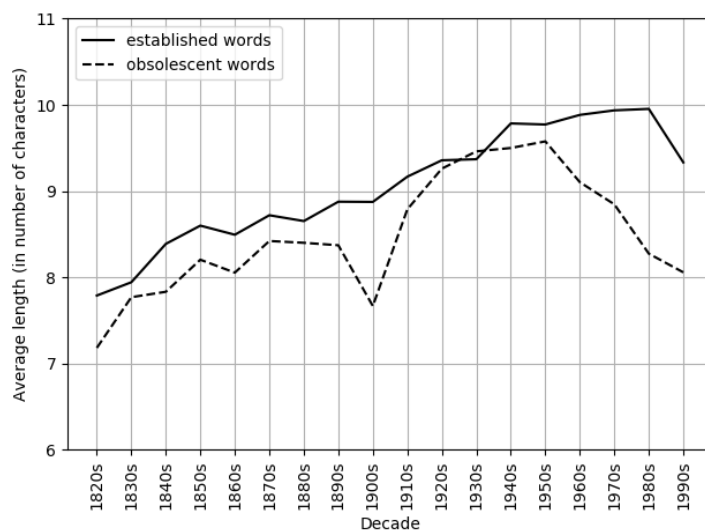


Figure 3.2: Average length (in number of characters) of words that became established and obsolete in given decades using a one-deviation criterion.

obsolete, the curve that represents nouns seems to exhibit a downward trend, while the others show constant fluctuation through time. Again, the fact that COHA is balanced by genre across time suggests that these patterns, in principle, should not be due to artifacts of the corpus, even though additional investigation is needed to better comprehend the phenomena reported here.

### 3.4.3 Case 3: Lexical heritage from past decades

The lexicon of every language at time  $t$  embodies strata from different periods in time during which new words that are still used at  $t$  became established. We are now onto a bit of “stratigraphy”, employing our algorithm to generate lists of those words established in different decades that are today the most popular in the corpus. More precisely, we select, from the words established in each

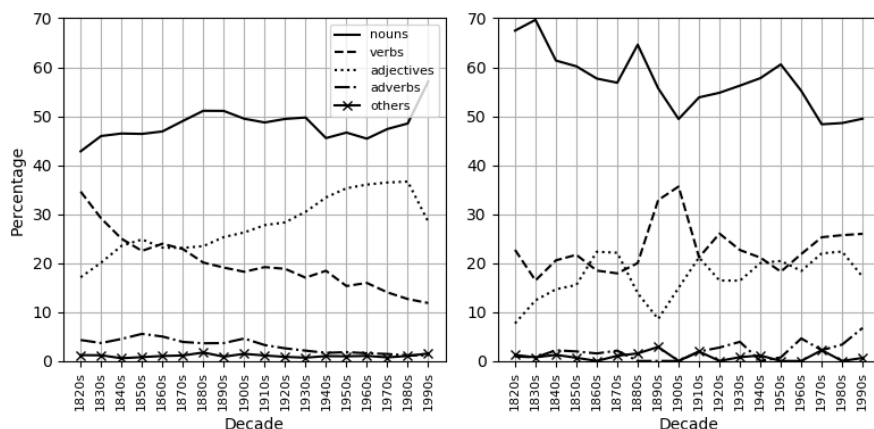


Figure 3.3: Percentage of different parts of speech among words that became established (left) and obsolete (right) in different decades.

decade between the 1850s and the 1980s, the fifty words that are most frequent in the 2000s. This ensures that we capture a portrait of today’s lexical heritage from past decades which is both reasonably detailed and still salient to speakers of American English.

The result of the selection procedure is displayed in Appendix A. After grouping these words into semantic categories (e.g. by using tools like Empath (Fast et al., 2016) or LIWC (Tausczik and Pennebaker, 2010)) or building networks (e.g. by a co-occurrence metric), it would be possible to make some generalizations concerning which semantic domains have been major contributors to these different historical strata or to determine the overall relationship among words established in a given decade<sup>10</sup>.

Impressionistically, for instance, it seems that the 1870s gave us

<sup>10</sup> These generalizations, however, should be made carefully. Even if the corpus is balanced by genre across time (which is the case of COHA), the topics covered by the texts themselves might vary systematically (and not nicely randomly) over time. A possible way of mitigating this (potential) issue could be to implement a topic detection method, such as latent Dirichlet allocation (LDA) (Blei et al., 2003), in order to ensure that topics are coherent over time.

much vocabulary relating to the built environment, such as *hallway*, *downtown*, *driveway*, *taxi*, *headlights* and *neon*. The 1880s were big on sports, cf. *golf*, *hockey*, *olympics*, *coaching*, *scoring*. The 1890s were innovative in the communication domain, see *movie*, *television*, *wireless*, *phones*. The 1910s opened the written language to include words that would have been considered too obscene to print earlier: *fuck*, *goddam*, *dick*. The 1920s introduced several relatively abstract concepts relating to workflow: *coordinator*, *feedback*, *processing*, *implementation*, *operational*. The list of the 2000s most frequent words stemming from the 1940s does not reveal that a big war happened; instead, for instance, we see elements of such an everyday affair as food consumption: *supermarket*, *microwave*, *fridge*, *burgers*, *yogurt*. The 1950s show nascent environmental concerns: *pesticides*, *recycling*, *environmentally*, *pollutants*. Apart from giving us the concept of *lifestyle(s)* in general, the 1970s also showed news in different domains of lifestyle, such as the food domain, cf. *tofu*, *fast-food*, *sushi*, *veggies*. During the 1980s, the (personal) computer is the single most dominant factor in lexical innovation: *laptop(s)*, *database(s)*, *pcs*, *algorithms*, *download*, *firewall*. As we move closer to the present, it is predictable that some of the lexical legacy encountered is going to be short-lived, and some of the terms from the 1970s and 1980s, indeed, already feel somewhat outdated.

A comprehensive study of the lexical legacy of different periods in current American English could probably use a larger selection and, as mentioned, systematic methods of defining semantic domains and networks. In any case, we show that our algorithm is effective for identifying the lexical material needed for such research.

#### 3.4.4 Case 4: Lost words

Another use of our proposed method is to generate lists of previously popular items that became obsolete in the corpus during

its time span. This is interesting because, unlike innovations, obsolete items are not commonly covered in existing literature (Tichý, 2018). Here, we select, from the words that became obsolete in each decade between the 1850s and the 1980s, the ten with the highest frequency before their obsolescence. In this fashion, we display a vocabulary that was particularly relevant in the past, but that has lost terrain in American English after some decades.

The lists of the once common words that became obsolete in COHA in a particular decade are shown in Tables 3.3 (1850s-1920s) and 3.4 (1930s-1980s). In the second half of the 19th century, it is possible to encounter typos and spelling mistakes that ceased to appear in the 20th century (maybe partially due to the development of more accurate typing and printing techniques), such as *had'nt* (1870s), *do'nt* (1880s), *hav'nt* (1890s), *was'nt* (1890s) and *did'nt* (1890s). There are also several other words that are still easily recognizable but have obsolete or semi-obsolete spellings, including *errour* (1850s), *pennyless* (1870s), *musquitoes* (1880s), *negociation* (1890s), *villany* (1930s), *reconnoissance* (1940s), *trowsers* (1950s) and *persistency* (1960s), to name a few. In some cases, the obsolete spelling is more faithful to the etymology of the word, as in *holydays* (1880s), which became “holidays”, and *cocoa-nut* (1930s), which became “coconut”. Further, the lists exhibit spellings that are still present in the corpus, but no longer with a specific syntactic function, such as *under* as a comparative adjective (1900s), *itself* as a singular common noun (1900s) and *notwithstanding* as a subordinating conjunction (1970s).

Of particular interest is the illustration provided by these lists of the phenomenon of the historical spelling change of English compounds. According to Shertzer (1996), “[t]he usual sequence is for the words to be written separate at first, then to become hyphenated, and finally to be written solid” (p. 109). We observe that several compounds that are nowadays usually written

Table 3.3: Lists of common words (+ PoS tags) in previous decades that became obsolete in the corpus in a particular decade (1850s-1920s). Words are ordered according to their frequency before their obsolescence. When the word ranked in the eleventh position has the same frequency as the one in the tenth position, we include both. The meaning of each PoS tag is explained in Appendix A.

1850s	1860s	1870s	1880s
errour_nn1	copy-right_nn1	had'nt_vv0	phrensy_nn1
scymetar_nn1	hazle_nn1	ancke_nn1	sassacus_nn1
almanzor_nn1	do'st_vv0	pennylesse_jj	anckes_nn2
pedrillo_nn1	pannels_nn2	wo-begone_jj	cotemporary_jj
musquetry_nn1	phrensied_jj	inartificial_jj	musquitoes_nn2
inquietudes_nn2	choaked_vvd	wrapp_nn1	afford_nn1
renegado_nn1	fire-side_jj	teaze_vvi	holydays_nn2
errours_nn2	barb'rous_jj	rivalships_nn2	gallopped_vvd
zegri_nn1	famish_jj	a'nt_vv0	apalachian_jj
broad-street_nn1	fann_nn1	returnless_jj	do'nt_vv0
potawatamies_nn2	incommunicative_jj		vanquish_jj
1890s	1900s	1910s	1920s
merchandize_vv0	shakspeare_vv0	shakspeare_nn1	the_nnt1
rivalship_nn1	immoveable_jj	eend_nn1	desponding_jj
hav'nt_vv0	under_jjr	did'st_vv0	flag-staff_nn1
had'st_vv0	say'st_vv0	creatur_nn1	sportively_rr
guarantied_vvn	xve_nn1	deth_vvz	befel_vv0
intenseness_nn1	itself_nn1	piano-forte_nn1	stopt_vv0
was'nt_vv0	wall-street_nn1	saidst_vv0	discomposed_vvn
negociation_nn1	pedee_nn1	thou'st_nn1	enginery_nn1
cretur_nn1	xve_vv0	applauses_nn2	school-fellows_nn2
did'nt_nn1	sdeath_nn1	knitting-work_nn1	sarvice_nn1
		see'st_vv0	

in a solid form are present in the corpus as hyphenated compounds and that these became obsolete at some point – probably around the time when their corresponding solid form were gaining popularity. This is the case of *copy-right* (1860s), *fire-side* (1860s), *wo-begone* (1870s) (now most commonly written “woe-begone”), *piano-forte* (1910s) (now mostly encountered as just “piano”), *flag-staff* (1920s), *school-fellows* (1920s), *dew-drops* (1930s),

Table 3.4: Lists of common words (+ PoS tags) in previous decades that became obsolete in the corpus in a particular decade (1930s-1980s). Words are ordered according to their frequency before their obsolescence. When the word ranked in the eleventh position has the same frequency as the one in the tenth position, we include both. The meaning of each PoS tag is explained in Appendix A.

1930s	1940s	1950s
villany_nn1	csar_nn1	trowsers_nn2
prison-house_nn1	new-comer_nn1	school-master_nn1
nuther_vv0	custom-house_nn1	mantel-piece_nn1
wofully_rr	custom-house_jj	despatch_vvi
unbiased_jj	bethink_vv0	hill-top_nn1
dew-drops_nn2	reconnaissance_nn1	aliment_nn1
cocoa-nut_jj	hill-tops_nn2	corner-stone_nn1
log-house_nn1	prayer-meetings_nn2	leipsic_nn1
can't_vv0	school-boys_nn2	exhaustless_jj
palm-tree_nn1	sketch-book_nn1	self-complacency_nn1
1960s	1970s	1980s
acquirements_nn2	sich_vv0	intrusted_vvn
inclosure_nn1	ball-room_nn1	arm-chair_nn1
persistence_nn1	now-a-days_rt	fellow-men_nn2
state-room_nn1	frying-pan_nn1	quitted_vvd
upon_nn1	notwithstanding_cs	with_nn1
intrenchments_nn2	hesitating_jj	common-place_jj
snuff-box_nn1	reprobation_nn1	fitly_rr
strifes_nn2	banditti_nn2	unwearied_jj
guard-house_nn1	by-gone_jj	small-pox_nn
heart-strings_nn2	plighted_jj	inclosed_vvn

*new-comer* (1940s), *corner-stone* (1950s), *state-room* (1960s), *ball-room* (1970s), *now-a-days* (1970s), *arm-chair* (1980s), *common-place* (1980s) and various others that can be recognized in the table. Nonetheless, a few compounds, such as *wall-street* (1900s) and *knitting-work* (1910s), seem to have taken the opposite direction, now being more commonly written as separate words. A comprehensive study aiming to analyze this phenomenon in a quantitative fashion could benefit from our proposed method to obtain these lists

of obsolete items per time frame and investigate how different factors (e.g. time, accumulated frequency, sudden frequency rise/fall) act and impact this process of orthographic variation and change.

### 3.4.5 Case 5: Short-lived words

The method described in Section 3.3.3 is able to assist in the identification of items classified as *established* or *obsolete*, but not of items evaluated as *short-lived*. Here, we provide a short case study in which we suggest a way of adapting it for this specific purpose. Our goal is to find words that flared up in the corpus for some time and then, still during the period covered by the corpus, disappeared. According to our previously mentioned criteria, these words are considered neither established (since they are already gone) nor obsolete (since they are not part of the corpus in its initial period), but in some cases it might be interesting to analyze them in order to investigate the process of lexical variation and change in more detail.

A possible way of adapting our method to the case of short-lived items is by applying the proposed algorithm to selected intermediate subcorpora. One solution would be to look for items whose diachronic sequences hold only 0s in their extreme time frames, such as in 0001111000, then cut off the extremes of the corpus (say, the  $n_1$  time frames in the beginning and the  $n_2$  time frames in the end of the time span covered by the corpus) and, finally, apply the algorithm only to the remaining intermediate sequences, looking for established, obsolete and permanent items in these subcorpora.

For the present exploratory purposes we adapted our method to handle cases of words that did not appear in COHA before the 1860s and disappeared again no later than the 1950s – in other words, these items are present neither in the five first nor in the five last time frames of the corpus ( $n_1 = n_2 = 5$ ). We then applied



our algorithm considering just this subsection of the corpus. We extracted words evaluated as permanent – which are, of course, perfect cases of short-lived words, presenting the diachronic sequence [00000]1111111111[00000]<sup>11</sup>. We also gathered other not-so-short-lived words evaluated as established and obsolete in the subsection of the corpus, but only those that appeared in at least eight decades and with no deviations allowed<sup>12</sup>.

The words that emerged from this analysis are listed alphabetically in Table 3.5. The vast majority of them are compounds (either hyphenated or solid), short-lived spelling variants and bona fide words that came and went. Among the hyphenated compounds, we find words such as *farm-lands*, *hair-pin* and *saddle-bag* – all of them more commonly written in a solid form nowadays. These data are useful for the study of the historical spelling change of English compounds mentioned in Section 3.4.4. Words such as *comp'ny*, *yisterday* and *s'posin* are examples of short-lived spelling variants. The comparative adjective *humaner* (meaning *more humane*) and the nouns *leisureliness* (*leisurely* + *-ness*) and *stereopticon* (an old type of slide projector) are interesting examples of short-lived items found here: when searching on another source, the Google Books Ngram Viewer<sup>13</sup>, we find that all of them exhibit a similar frequency pattern, peaking around the 1920s.

These results are just an illustration of the kind of content that can be obtained from such an analysis. It is important to notice that looking for short-lived items is not, in principle, one of the goals of the method introduced in this chapter, and that the adaptation presented in this case study is just a workaround. The main pitfall of this adaptation is that it depends on the selection of spe-

<sup>11</sup> The 0s in between square brackets correspond to the extremes of the corpus that were cut off.

<sup>12</sup> That is, those which, for the period of the subcorpus studied, presented the diachronic sequences 0111111111, 1111111110, 0011111111 and 1111111100.

<sup>13</sup> <https://books.google.com/ngrams/> .

Table 3.5: Words (+ PoS tags) classified as short-lived according to the adaptation of our method and considering the period between the 1860s and the 1950s. Words are alphabetically ordered. The meaning of each PoS tag is explained in Appendix A.

a-beatin_nn1	crep_nn1	ha'r_nn1	race-track_nn1
a-laughin_nn1	dilapidated-looking_jj	hair-pin_nn1	rose-petals_nn2
a-puttin_nn1	dish-towels_nn2	hay-wagon_nn1	s'posin_nn1
a-quiver_vv0	dust-heap_nn1	hereinbefore_rr	sabe_vvi
a-sittin_nn1	ear-drums_nn2	herse'f_nn1	saddle-bag_nn1
all-rail_jj	earnin_nn1	hez_vv0	spoilin_nn1
alongshore_nn1	east-bound_jj	high-tariff_jj	staff-officer_nn1
baggage-man_nn1	farm-hands_nn2	humaner_jjr	station-master_nn1
bath-chair_nn1	farm-lands_nn2	ice-floe_nn1	stereopticon_nn1
bird-shot_nn1	field-glass_nn1	idealizing_jj	street-cars_nn2
black-fringed_jj	field-glasses_nn2	jumping-jack_nn1	talesmen_nn2
bodder_vvi	fitten_vvn	leisureliness_nn1	tek_vvi
bofe_nn1	food-supply_nn1	lucile_nn1	trades-union_nn1
bread-winner_nn1	foregathered_vvd	myse'f_nn1	unfoldment_nn1
broncho_nn1	forehanded_vvn	pack-train_nn1	up-train_nn1
burled_vvn	four-bit_jj	pay-rolls_nn2	w'at_nn1
catchee_nn1	full-armed_nn1	pepsin_nn1	w'en_jj
chromos_nn2	garden-party_nn1	play-actin_nn1	water-bottle_nn1
coat-sleeves_nn2	glarin_nn1	pony-cart_nn1	weazened_vvd
comp'ny_jj	groceryman_nn1	prohibitionist_jj	wedding-bells_nn2
consul-general_jj	grouped_jj	pulse-beats_nn2	yisterday_nn1

cific subsections of the corpus to be analyzed by the researcher. A possible goal for future work is to design and develop a specific and more effective method for finding short-lived items in diachronic corpora.

### 3.5 Concluding remarks

In the field of corpus linguistics, the analysis of diachronic corpora with the goal of explaining diverse phenomena in human languages is becoming increasingly widespread. In this context, we need methods and procedures aiming to discover trends and patterns in the dynamics of a language as we process big amounts of text com-

putationally. With the present contribution, we hope to specifically generate more interest in the birth and death of components such as words, expressions and grammatical constructions in corpora that span over time.

Here, we introduce the notions of *establishment* and *obsolescence* as complementary to the trivial concepts of first and last attestations of linguistic items in diachronic corpora. Subsequently, we propose an algorithm to identify the time period of establishment and obsolescence of linguistic items based on their frequency in a diachronic corpus. This algorithm may be employed for the analysis of any linguistic item, be it lexical, phonological or morphosyntactical. The method proposed here is, of course, only one of the numerous possibilities for the achievement of similar goals. Other methods, including more mathematically sophisticated ones, could be evaluated as well. Alternatives that look promising for further consideration are approaches that would model *probabilities* of establishment and obsolescence. Such approaches would have the double advantage of allowing more accurate estimates of when a probability of occurrence exceeds a given threshold, and of allowing to make such estimates with fewer arbitrary parameters (e.g. lengths of periods, occurrence thresholds within a period, which patterns to consider as indicating what kind of event etc.). In this work, our focus is to demonstrate a simpler and easier-to-implement method, but we plan to discuss more sophisticated approaches in future studies.

We demonstrate the applicability of our proposed algorithm using a real corpus spanning 200 years of data and supplying case studies concerning the character of words that got established and obsolete in American English in different periods. Among the outcomes of these case studies is the observation that the percentage of established words among all words across decades fluctuates without showing a specific upward or downward trend. We also found

that the proportion of adjectives among new words has increased steadily over the past two centuries, mostly mirrored by a decrease in the proportion of new verbs. Then, we provided a sketch study of the lexical heritage in American English, identifying words that became established in different decades and are still frequent in the 2000s. We also looked at obsolescent vocabulary – vocabulary that was previously frequent but has been getting lost over the decades. Finally, we briefly investigated whether the method could be adapted to find short-lived words – words that flared up in the corpus for some time and then disappeared. These sketch studies are mainly presented with the goal of motivating future studies employing the method presented here.

It may be obvious but still it is necessary to recall that a corpus is different from a language. As a consequence, when we consider the establishment or the obsolescence of a linguistic item in a *corpus*, we are not necessarily referring to the establishment or the obsolescence of this item in a *language*. This distinction is particularly relevant when we deal with corpora based on written texts (like COHA itself or the Google Books corpus) – since, for instance, an item might be used for a long time in the oral language before it gets established in the written register. When considering the whole language, it is clear that the algorithm can only identify the decade during or *before which* (*ante quem*) a word became established or the decade during or *after which* (*post quem*) a word became obsolete. This situation is of course due to the fact that “it is much simpler to prove that something exists (...) than to prove that something does not exist” (Tichý, 2018, p. 82). This fact becomes even clearer if we think about the application of our method to domain-specific corpora (consisting of academic, legal, medical etc. texts): the results will of course reflect the specificity of the analyzed data.

Regarding our case studies, it is important to remember that words are pairings of form and function. Words not always start

their lives with a meaning and get lost with that same meaning, since in real-life diachronic lexical change there are also forms that come into being with a particular connotation but at some point lose that connotation, while still living on with a completely different one; and such words occur alongside words that live on with their original meaning. This must be taken into account when the researcher employs our (or any other) method to automatically obtain lists of forms that get established or become obsolete in a corpus.

As stated by Hilpert and Mair (2015), it is imperative to demonstrate “how the use of corpus data allows researchers to go beyond the mere statement that a grammatical change happened, and to address the questions of **when** and **how** something happened” (Hilpert and Mair, 2015, p. 199, emphases in original). With our theoretical discussion, our proposed algorithm, and the case studies that were presented here, we hope to have taken a step in this direction.

