

Analogy at the level of phonology: the emergence of intrusive-*r* in English

The main goal of this presentation is to show that a memory-based analogical model can help us understand the emergence of intrusive-*r* in Southern British English. The apparent unnaturalness and cross-linguistic rarity of intrusive-*r* has led many researchers to conclude that it is synchronically arbitrary and that its present behaviour can only be fully understood in the context of its historical development (McCarthy, 1991; Blevins, 1997; Halle & Idsardi, 1997; McMahon, 2000; Gick, 2002). To be more specific, most of the authors cited above claim that the present situation is the result of rule-inversion: an original rule of *r*-deletion gave rise to a rule of *r*-insertion. However, it is not clear what the precise conditions are under which this type of rule-inversion can take place and why it does not occur in many other languages with comparable deletion rules. I argue that an analogical approach which takes frequency into account can lead to a better understanding of this unnatural development and present a number of simulations which support this view.

The basic claim is that the emergence of intrusive-*r* is the result of the analogical extension of the alternating pattern of behaviour exhibited by etymologically *r*-final words (henceforth the *r*-ful class; e.g. ‘bar’: [ba:] in C and # vs. [ba:r] in V) to the class of words ending in [ə], [ɑ:] and [ɔ:] (henceforth the *r*-less class; e.g. ‘idea’ [aɪ'diə], ‘ma’ [mɑ:], ‘raw’ [rɔ:]) based on phonetic similarities between these two groups (both *r*-ful and *r*-less forms end in non-high vowels in preconsonantal and prepausal position). However, an analysis relying solely on the partial merger between the *r*-ful and the *r*-less class cannot explain the direction of the extension, and also leaves open the question of why we do not find similar changes in other languages with similar deletion rules. Therefore, it will be useful to take a closer look at the frequencies within these two classes (1) and also the frequencies of forms in different environments within the *r*-less group (2) (the token frequencies below are taken from a 2 million word phonetically transcribed corpus of 18th century English—the CE18 corpus—compiled by the author):

(1)				(2)			
	R-LESS	R-FUL	RATIO		_V	_C	_
ə#	1,553	99,881	1:64.31	ə#	422	822	274
ɔ:#	1,487	51,871	1:34.88	ɔ:#	421	971	86
ɑ:#	112	9,397	1:83.90	ɑ:#	19	63	18
SUM	3,152	161,149	1:51.13				

The overwhelming frequency difference between the *r*-ful and the *r*-less classes and the relative infrequency of prevocalic forms within the *r*-less class provide a straightforward explanation for the emergence of intrusive-*r*: as the speakers have very little evidence to support any hypothesis about the prevocalic behaviour of etymologically *r*-less words, they will often attribute an *r*-ful pattern to them (which can be regarded as a default pattern among words whose citation form ends in a non-high vowel). It should be clear that this type of extension should only occur in languages where these particular frequency relations hold and that it cannot proceed in the opposite direction.

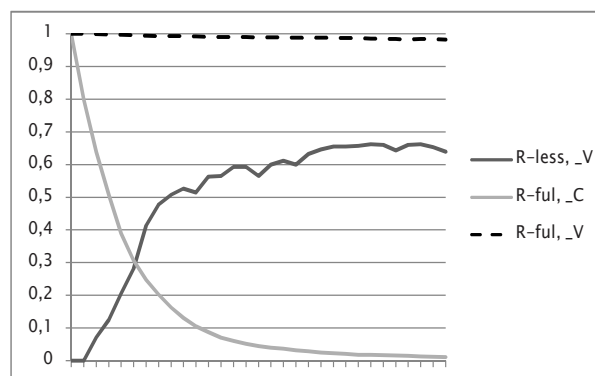
While the argument above seems to be in line with what we know about analogical change (Bybee 2001), it is formally rather inexplicit and therefore difficult to evaluate. It would be interesting to see whether any existing formally explicit analogical model would predict this type of extension to occur under the conditions outlined above. To test this, I have compiled a dataset consisting of the last five sounds of the citation forms of all the types in the CE18 corpus along with their patterns of behaviour (e.g. [=, =, k, æ, t] {no alternation}, [a, I, d, I, @] {no alternation}, [=, b, E, t, @] {*r*~zero}) and used a memory-

based learner, TiMBL (Daelemans et al., 2007), to predict the behaviour of each word in the dataset on the basis of all the other forms. The results were as expected: TiMBL correctly predicted a non-alternating pattern for words ending in a high vowel or a consonant other than *r* (in 100% of the cases) and an *r*~zero pattern for *r*-ful forms (in 99% of the cases); however, TiMBL incorrectly predicted an *r*~zero pattern for *r*-less forms ending in a non-high vowel (in 100% of the cases), that is, it extended the *r*-ful pattern to this class of words. This means that the extension described in the previous paragraph can be replicated by formally explicit analogical models.

However, there are two major problems with the simulation above. First, it predicts a sharp transition from a pre-intrusion dialect to one with an exceptionless insertion rule. This is not supported by the historical record, which shows that intrusive-*r* emerged gradually as rhoticity declined (Hay & Sudbury, 2005). Moreover, intrusive-*r* is clearly not categorical in modern varieties English (Foulkes 1998). The second problem stems from the fact that TiMBL was designed with simple categorisation tasks in mind, which means that (i) the input dataset has to be specified in terms of types rather than tokens of use and (ii) the pattern of behaviour characteristic of each type has to be specified explicitly. This raises several issues. Types are abstractions over sets of tokens, which means that they cannot be associated with a single phonetic form—I could have chosen the prevocalic forms of the types in the CE18 corpus to represent them in the simulation, in which case there would have been no extension (as the *r*-ful and *r*-less classes are fully distinct prevocalically). Moreover, each type has to be specified as belonging to a single category (i.e. alternating or non-alternating), which results in the loss of all information about word-specific patterns of variation—although there are several studies indicating that such patterns exist in the case of intrusive-*r* (e.g. Hay & MacLagan to appear). Finally, by using explicit behavioural patterns we reintroduce generative rewrite rules into the model through the backdoor—not necessarily a problem in itself, but clearly incompatible with the basic claims of most memory-based models (e.g. Bybee 2001).

The solution to these problems is to construct a dataset consisting of tokens of use rather than types and use an algorithm that can extract patterns of behaviour by looking at semantic and phonetic relations between the tokens themselves. The model used in this presentation is a combination of four-part analogy (e.g. Lepage 1998) and the exemplar-theoretic framework presented in Nosofsky (1986). The dataset consisted of 1 million words chosen from the CE18 corpus; the algorithm went through all the tokens within the dataset and tried to find a suitable phonetic form for each of them on the basis of their phonetic environment ($__C$, $__V$, $__#$), a semantically identical form in a different environment and a phonetically similar analogical model. For instance, when the model had to produce a prevocalic token of ‘idea’, it took a preconsonantal token of the same word ([aɪdɪə]), looked for a phonetically similar form in preconsonantal position (the analogical model; e.g. [dɪə]) and a token of the analogical model in prevocalic position (e.g. [dɪər]), and applied the difference between the two forms of the analogical model to the preconsonantal token of ‘idea’, to finally output [aɪdɪər]. This is shown in (3) below. This resulted in a small number of errors (i.e. extensions of the *r*-ful pattern), which were consolidated into more robust patterns through repeating the simulation several times, always taking the output of the previous simulation as the input of the next one. The simulation started with a fully rhotic dialect, with a 20% bias for word-final *r* to be lost preconsonantly and prepausally and eventually produced a dialect with both linking-*r* and intrusive-*r* (4). Thus, a token-based analogical model can simulate the emergence of intrusion in SBE in a realistic way.

(3) 'dear' [dɪə] → [dɪər] (4)
 'idea' [ajdɪə] → [ajdɪər]



References

- Blevins, J. (1997). Rules in Optimality Theory: Two Case Studies. Roca, I. (ed.), *Derivations and constraints in phonology*, Oxford University Press, Oxford, pp. 227–260.
- Bybee, J. L. (2001). *Phonology and language use*. Cambridge University Press, Cambridge.
- Daelemans, W., J. Zavrel, K. van der Sloot & A. van den Bosch (2007). TiMBL: Tilburg Memory Based Learner, version 6.1, Reference Guide. ILK Research Group Technical Report Series no. 07-07.
- Foulkes, P. (1998). English [r]-sandhi: a sociolinguistic perspective. *Leeds Working Papers in Linguistics & Phonetics* 6, pp. 18–39.
- Gick, B. (2002). The American intrusive l. *American Speech* 77, pp. 167–183.
- Halle, M. & W. J. Idsardi (1997). r, hypercorrection and the elsewhere condition. Roca, I. (ed.), *Derivations and constraints in phonology*, Oxford University Press, Oxford, pp. 331–348.
- Hay, J. & M. MacLagan (to appear). Social and phonetic conditioners on the frequency and degree of ‘intrusive /r/’ in New Zealand English. Preston, D. & N. Niedzielski (eds.), *A sociophonetics reader*, Mouton de Gruyter, Berlin.
- Hay, J. & A. Sudbury (2005). How rhoticity became /r/-sandhi? *Language* 81, pp. 799–823.
- Lepage, Y. (1998). Solving analogies on words: an algorithm. *Proceedings of COLING-ACL’98, Montreal, August 1998, vol. I*, pp. 728-735.
- McCarthy, J. J. (1991). Synchronic rule inversion. Sutton, L., C. Johnson & R. Shields (eds.), *Proceedings of the Seventeenth Annual Meeting of the Berkeley Linguistics Society*, Berkeley Linguistics Society, Berkeley, CA., pp. 192–207.
- McMahon, A. (2000). *Lexical Phonology and the History of English*. Cambridge University Press, Cambridge, UK.
- Nosofsky, R. M. (1986). Attention, similarity and the identification-categorization relationship. *Journal of Experimental Psychology: General* 115, pp. 39–57.