

Probabilities improve stress-prediction in a CFG of Hawaiian phonology

‘Ōiwi Parker Jones

Phonetics Laboratory

41 Wellington Square

Oxford, OX1 2JF

United Kingdom

oiwi.parkerjones@ling-phil.ox.ac.uk

1 Introduction

Should probability play a role in linguistics? Where Chomsky (1957: 16–17) influentially rejected probability in syntactic theory, “[i]n phonological theory, probability has not so much been rejected as disregarded” (Coleman, 2003: 89). Disregard has, however, given way to a growing literature on the use of probability across the various linguistic sub-disciplines (see, e.g., Bod et al., 2003; Coleman, 2003).

This paper is a case-study of probability in phonology, both as it applies to an improved description of Hawaiian stress-assignment, and as this description, in turn, reflects back on the probability question, above.

2 Grammars

By formalizing two strongly equivalent analyses, where one is a non-probabilistic Context-Free Grammar (CFG) and the other is a Stochastic Context-Free Grammar (SCFG) (Booth, 1969; Suppes, 1970), we can put the probability question to a test. For a given data set, if the SCFG does not outperform its strongly equivalent CFG, then parsimony should compel us to reject, rather than disregard, the added probabilities.

On the other hand, should the SCFG outperform its strongly equivalent, non-probabilistic CFG, then we ought, at least, to accept some role for probability in phonology; this would support the growing literature mentioned above.

Let our data set be Hawaiian. Schütz (1978, 1981) argues that Hawaiian stress-assignment is not 100% predictable, based on words like /ma.ku.a.'hi.ne/ ‘mother’ and /,ʔe.le.ma.'ku.le/ ‘old man’. It might help to illustrate this argument by developing Schütz’s analysis into a CFG.

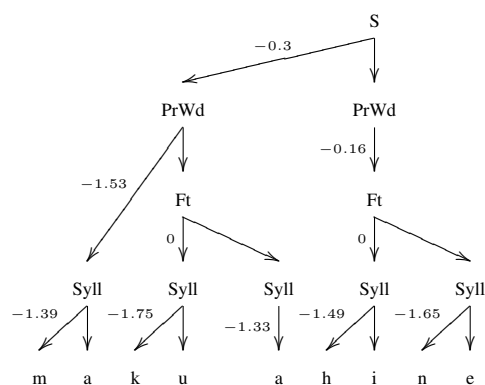


Figure 1: /ma.ku.a.'hi.ne/ parse-tree.

The crucial parse-trees for ‘mother’ and ‘old man’ are in Figures 1–4. Note that the terminal symbols are phonemes. The non-terminal symbols are syllables (Syll), metrical-feet (Ft), prosodic words (PrWd), and the start symbol (S). Also note that the leftmost syllable in each metrical-foot is stressed. The rightmost stress in a word is primary. Finally, let us ignore the labeled branches for the moment, as they do not apply to the non-probabilistic CFG.

The argument follows. The correct parse for ‘mother’ (Figure 1) is paralleled by an incorrect parse for ‘old man’ (Figure 2); except for their terminal expansions, these parse-trees have the same structure. Thus, the correct parse for ‘mother’ implies an incorrect parse for ‘old man’. Moreover, the correct parse for ‘old man’ (Figure 3) implies an incorrect parse for ‘mother’ (Figure 4). Therefore, no matter how we order the CFG rules, a procedural interpretation will get either ‘mother’ or ‘old man’ wrong. Hence, Hawaiian stress-assignment is not 100% predictable.

Now, this conclusion might be true; as we shall

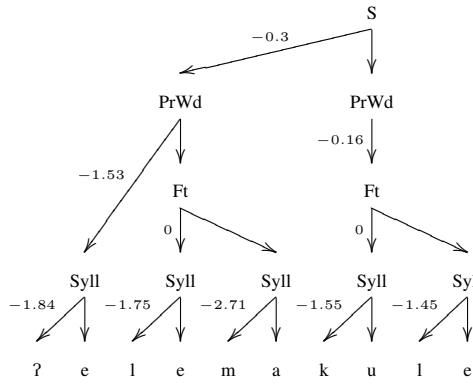


Figure 2: */ʔe.le.ma.'ku.le/ parse-tree.

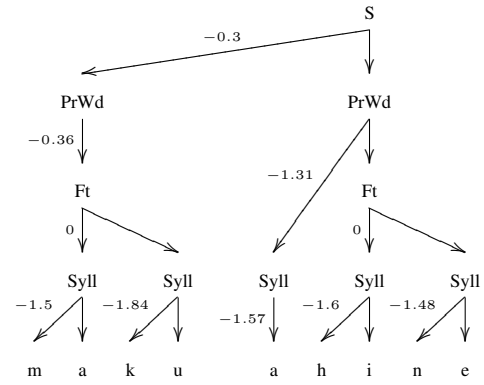


Figure 4: */ma.ku.a.'hi.ne/ parse-tree.

see, neither grammar gets 100% of the evaluation data right. But the SCFG still performs better than its strongly equivalent CFG at predicting stress-assignment in Hawaiian.

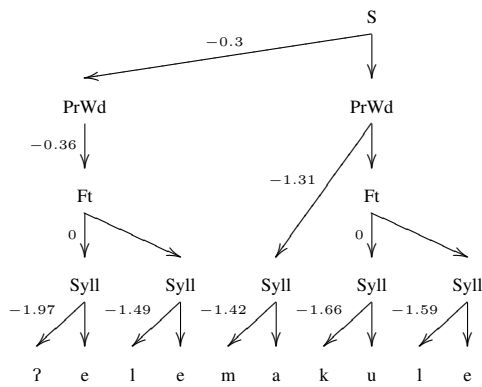


Figure 3: /ʔe.le.ma.'ku.le/ parse-tree.

Returning to Figures 1–4, each labeled branch expresses the base-10 log probability for some SCFG rule, where the probabilities were obtained by training the grammar on data from a Hawaiian dictionary (Pūku'i and Elbert, 1986). The probability of a parse-tree is just the sum probability of its rules, so Figure 2's probability is -11.29 . By contrast, Figure 3's probability is -10.09 . The SCFG correctly picks /ʔe.le.ma.'ku.le/ over */ʔe.le.ma.'ku.le/, since a log probability of -10.09 is higher than a log probability of -11.29 . Moreover, the SCFG correctly picks /ma.ku.a.'hi.ne/ over */ma.ku.a.'hi.ne/, since a log probability of -9.59 is higher than that of -9.95 . In both examples, the SCFG correctly disambiguates the parses.

3 Evaluation

In a computational evaluation of 16,900 Hawaiian words, the CFG correctly parsed 84.6%. However, the SCFG correctly parsed 97.38%. These results demonstrate that probabilities improve stress-prediction in a CFG of Hawaiian phonology. There is a role for probability in phonology.

Acknowledgments

Thanks to John Coleman for supervision, Keola Donaghy for the electronic Hawaiian dictionary, Greg Kochanski for discussions on probability, and Lamakū for financial support.

References

- R. Bod, J. Hay and S. Jannedy, eds. 2003. *Probabilistic Linguistics*. MIT Press, Cambridge, MA.
- T. L. Booth. 1969. Probabilistic representation of formal languages. *Tenth Annual IEEE Symposium on Switching and Automata Theory*. Pp. 74–81.
- N. Chomsky. 1957. *Syntactic Structures*. Mouton, The Hague.
- J. S. Coleman. 2003. Commentary: probability, detail and experience. In J. Local, R. Ogden, and R. Temple, eds. *Phonetic Interpretation: Papers in Laboratory Phonology VI*. Cambridge University Press, Cambridge. Pp. 88–100.
- M. K. Pūku'i and S. H. Elbert. 1986. *Hawaiian Dictionary*. University of Hawai'i Press, Honolulu.
- A. J. Schütz. 1978. Accent in two Oceanic languages. *Anthropological Linguistics*, 20(4): 141–149.
- A. J. Schütz. 1981. A reanalysis of the Hawaiian vowel system. *Oceanic Linguistics*, 20(1): 1–43.
- P. Suppes. 1970. Probabilistic grammars for natural languages. *Synthese*, 22: 95–116.