# Machine Learning in Prosody and Meaning

*Shirley Fong, Madeleine Guettler, Sofia Hirschmann*
*Simmons University, Boston, MA, 02115*

## Abstract

Prosody in spoken language is widely thought to convey meaning, but the complexity and variety of prosodic renditions have made modeling a specific prosody-meaning mapping elusive. In this work, 238 prosodically labeled utterances are used to classify speech as either exclamative or neutral. A first pass was made to cluster accent types (emphasized words) into three groups. This research project captures the importance of a variable number of loosely categorizable elements; clustering is performed; and the percentage of accent type is used as an attribute in the final random forest classification step.

## Methods

The raw acoustics of pitch accents matter but are quite varied. In addition, there are more than one pitch accent per sentence and the ultimate goal is to classify the type of sentence. As a result, we seek to obtain a concise subset of pitch accent types that we can use in subsequent classification.

**Process:**
- Created a data frame with just the numerical attributes
- Scaled the data
- First trial: did not remove outliers
- Second trial: removed three outliers
- Tested important features by dropping features one at a time, while isolating features we thought would be the best performing and least redundant
- Calculated the sum of squares for each cluster and set of attributes
- Obtained the centroid/means for the data points for the best clusters
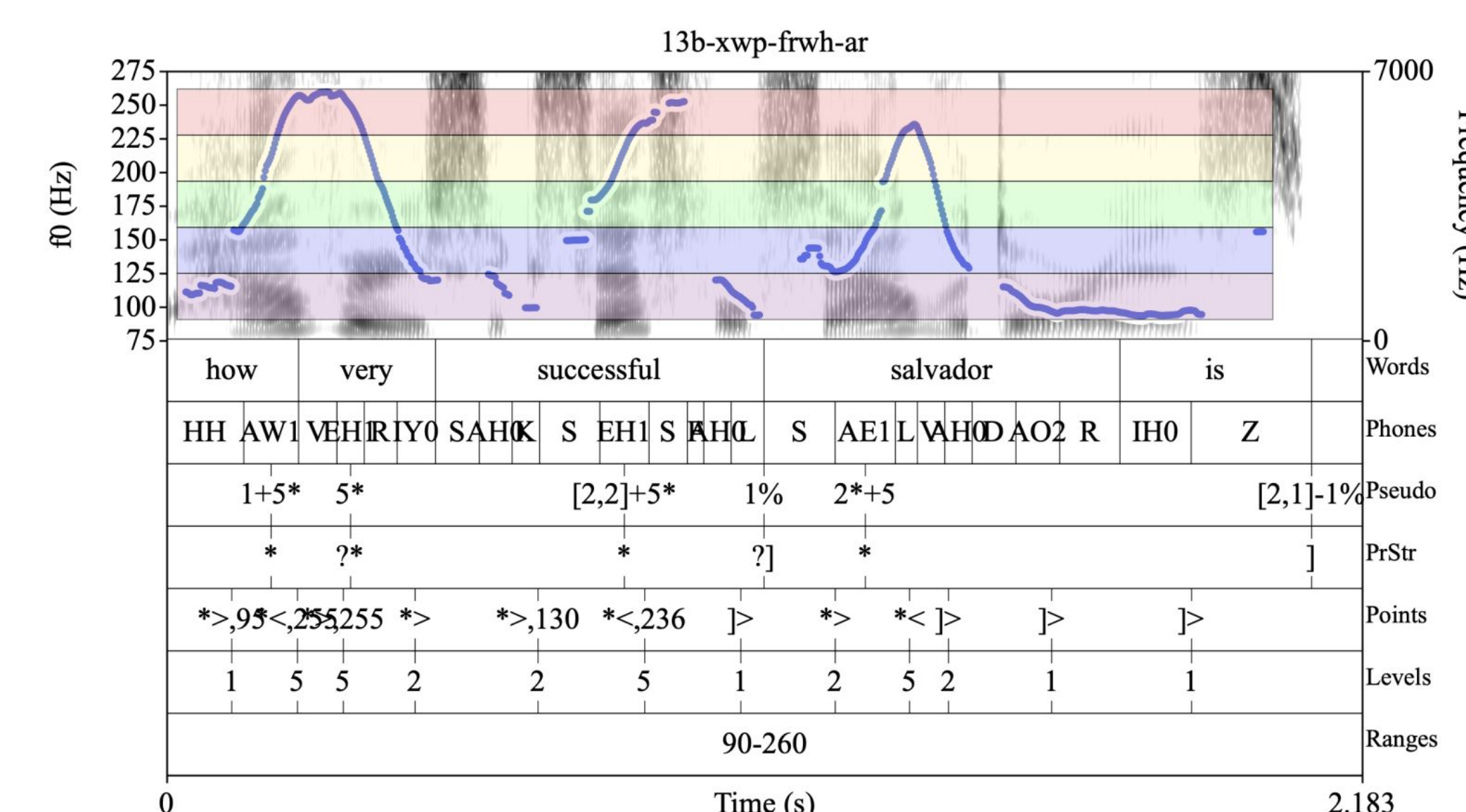
## References

**Selected References:**
- Beckman & Hirschberg. 1994. The ToBI annotation conventions.
- Ahn et al. 2021. PoLaR Annotation Guidelines (version 1.0). Available at https://osf.io/usbx5.
- Rett & Sturman. 2021. Prosodically marked mirativity. In Proceedings of WCCFL 37.
- Barnes, Veilleux, Brugos, & Shattuck-Hufnagel. 2012. "Tonal Center of Gravity: A global approach to tonal implementation in a level-based intonational phonology." *Laboratory Phonology* **3(2)**, pp. 337-383.

## Background



### What is prosody?
- In spoken language it's thought to convey meaning
- It's not what you say, but *how* you say it through alterations in pitch, duration, and intensity
- Prosody maps to meaning, and here our meaning is what we call mirativity: the idea of being surprised

### Question
- What are the best attributes and optimal amount of clusters for our k-means model to predict mirativity?

## Results



- Linguistics have long postulated several categorically different pitch accents and describe them in terms of the F0 (High or Low) and alignment (* placed early, medially or late with respect to the accented syllable)
- The Tonal Center of Gravity is a global measurement that describes alignment in time (TCoG-T) and F0 )TCoG-F)
- In the k=3, there appears to be a late and low accent that would correspond to an L*+H and two medially aligned accents: one relatively low in pitch (L+H*) and the other one relatively high (H*)
- As we increase the clusters, the clusters seem to separate variants of these three well known accent types