



BOOSTING OF PROSODIC AND PRONUNCIATION FEATURES TO DETECT MISPRONUNCIATIONS OF NON-NATIVE CHILDREN



Christian Hacker¹, Tobias Cincarek², Andreas Maier¹, André Heßler¹, Elmar Nöth¹

¹Institute for Pattern Recognition (LME), University of Erlangen-Nuremberg, Germany

²Graduate School of Information Science, NAIST, Nara, Ikoma, Japan

Outline

- German children reading English
- Acoustic models: trained on the Pf-Star native corpus (children from Birmingham)
- Pronunciation scoring with **176 features**
- Features from **Prosody/Pronfex** Module
- Feature selection with AdaBoost

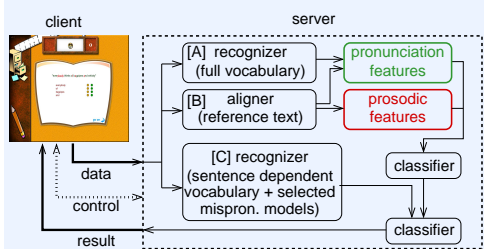


Figure 1: Caller: Computer assisted language learning from Erlangen. Here: Focus on paths via [A] and [B]

Non-Native Data

- 28 children, age 10–11 (English as 2nd lang. since 6 months)
- 72 min; vocabulary: 8088 tokens, 942 types
- reading errors, repetitions, word fragments, non-verbals

10 Raters:

- 1 graduate univ. student of English
- 8 German teachers of English marked words where they would have stopped the student in school
- 1 native British teacher (N)
- strictness of raters: 4.4% – 5.2%; N: 7.6% marked as errors

Reference:

Marked if at least 3 experts agree (5.6%, cf. strictness)

Class-wise-averaged recognition rate:

$CL = 0.5(REC_w + REC_c)$ (= av. Recall)
 c : correctly pronounced; w : wrongly pronounced

Agreement:

- rater vs. reference : $CL = 77\%$ (open average)
- pairs of raters: inter-rater: $CL = 70\%$; intra-rater.: 78 – 80%
- $REC_w \leq 78\%$; $REC_c \leq 99\%$

AdaBoost

- Select **weak classifiers** $h_t(\cdot)$ that use complementary info
- Here: each $h_t(\cdot)$ is **trained on exactly 1 feature**
- $h_t(x) = 1$ if mispronounced; 0 else

- Optimal threshold for each $h_t(\cdot)$ (criterion: CL)
- A weight $w_{0,i}$ is assigned to each word i of the training data. Weights of either class are distributed uniformly and sum up to 0.5.
- Choose the weak classifier $h_t(\cdot)$ with **lowest error** ϵ_t :
Words that are wrongly classified contribute with $w_{t,i}$ to the error.
- Use greater weights for all wrongly classified words:

$$w_{t+1,i} = w_{t,i} \frac{1 - \epsilon_t}{\epsilon_t}; \quad \alpha_t = \log\left(\frac{1 - \epsilon_t}{\epsilon_t}\right) \quad (1)$$

- Normalize the weights; $t = t + 1$; goto 3.
- Finally, combination to a **strong classifier**:

$$x \text{ is mispronounced, if } \sum_t \alpha_t h_t(x) \geq \frac{1}{2} \sum_t \alpha_t \quad (2)$$

- Leave-one-speaker-out** (loo) evaluation
- Calculate **mean** α_t for each feature over 28 loo-iterations to get this ranking:

Feature	Mean per word: $1/N \sum_{j=1}^N c_j$
Forced alignment	Word score
Speech recognizer	Posterior score [-1,0]
Phone bigram prob. of recognized phone sequence	Posterior score of reference word
Phone confusion $c_j = P(q_j p_j, M_w) / P(q_j p_j, M_c)$	FFT coefficient 1
q_j recognized phone	Deviation of phones (scatter)
p_j phone in reference	Posterior score of reference word
$j = 1 \dots N$: Index of phone	FFT coefficient 0
	Regression of the f_0 [-1,0]
	Bigram prob. of phones / #phones
	Position of the max. f_0 [1,1]
	Minimum [-2,-1]
	Normalized [-1,0]
	Maximum per word: $\max_j c_j$
	Mean [1,2]
	Minimum per word: $\min_j c_j$

Figure 2: Top 15 features selected with AdaBoost

Features Extraction

- Recognized word chain**, cf. Fig. 1 [A]
 - native models, LME-recognizer, 46% word acc.
 - 2500 additional mispronunciation models
- Forced alignment**, cf. Fig. 1 [B]
- Duration and Energy statistics**
 - estimated on native data
- Phoneme bigram model**
 - estimated on reference texts plus further data
- Phone confusion statistics on**
 - mispronounced words M_w
 - correctly pronounced words M_c

The Pronfex Module (63 features)

- Rate-of-speech, long pauses
- Duration (deviation from native statistics)
- Log-likelihood of reference word
- Word/phone accuracy and correctness
- Confidence of reference (N-best lists)
- Phone bigram probability
- Phone confusion using M_w and M_c

The Prosody Module (113 features)

- Pitch (F0), energy, duration, jitter, shimmer, pauses (e.g. min, max, mean, regression)
- Position of maxima, minima, on-set, ...
- FFT coefficients of the energy
- Context: preceding word [-1,-1], 2 succ. words [1,2]

Results

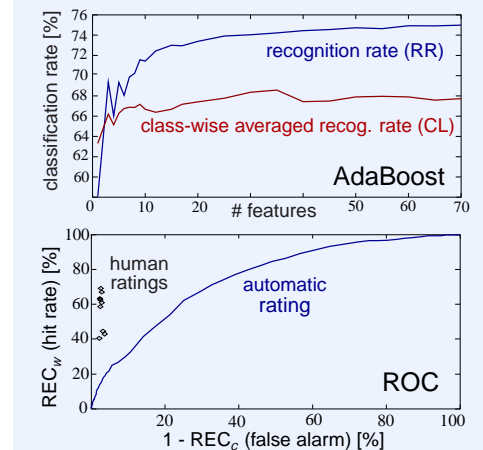


Figure 3: Classification rate with different numbers of features; ROC evaluation with 35 features

- Low CL with prosodic features, but useful extension to pronunciation features
- Best feature: **phone confusion**
- AdaBoost: **No overfitting** to training data
- Similar feature sets are selected
 - in all loo-iterations
 - using different references/experts
- 15 features: 66.7% CL
35 features: **68.6% CL**
→ **89% of human expert agreement**
- Teachers have high agreement on correct words (low hit rate on mispron.'s)