

Pronunciation Assessment Based upon the Compatibility between a Learner's Pronunciation Structure and the Target Language's Lexical Structure

Nobuaki MINEMATSU

Graduate School of Information Science and Technology, University of Tokyo

mine@gavo.t.u-tokyo.ac.jp

Abstract

Native-sounding vs. intelligible. This has been a controversial issue for a long time in language learning and many teachers claim that the intelligible pronunciation should be the goal. What is the physical definition of the intelligibility? The current work shows a very good candidate answer to this question. The author proposed a new paradigm of observing speech acoustics based upon structural phonology, where all the kinds of speech events are viewed as an entire structure and this structure was shown to be mathematically invariant with any static non-linguistic features such as age, gender, size, shape, microphone, room, line, and so on. This acoustic structure is purely linguistic and the phoneme-level structure is regarded as the pronunciation structure of individual students. This structure is matched with another linguistic structure, the lexical structure of the target language, and degree of compatibility between the two different levels of structures is calculated, which is defined as the intelligibility in this work. To increase the intelligibility, different instructions should be prepared for different students because no two students are the same. The proposed method can show the order of phonemes to learn, which is appropriate to a student and different from that of the others.

1. Introduction

There exist many kinds of English pronunciations socially accepted as intelligible all over the world, although some of them are clearly different from the native pronunciations. Many teachers claim that the intelligible pronunciation should be the goal of pronunciation training because pronunciation is just a tool for smooth speech communication. But it is very difficult to define the intelligible pronunciation physically because the intelligibility depends upon listeners. Especially in the case of non-native listeners, it is highly expected that different mother tongues will define different intelligible pronunciations. In spite of this difficulty, some bold attempts were made to discuss the intelligible pronunciation [1, 2], where non-native utterances were directly presented to listeners who were asked to repeat or type what they heard. A large number of facts of miscommunications were observed and, based upon the facts, the intelligible pronunciation was discussed. According to [2], acoustic and linguistic analysis of the facts implied that the most influential factor on the intelligibility is speech rhythm involved in an utterance. In both works, the listeners were Americans, which are just one candidate of the listeners, and this approach may have to continue until everybody on earth joins the experiment.

In this paper, another approach is taken, where the intelligibility is defined quantitatively with no attention to the listeners. The author proposed a new paradigm of observing speech acoustics based upon phonology [3]. Speech events are modeled

probabilistically as distributions, distance between any two of the events is calculated based upon information theory, and the events are relatively captured as a structure. The resulting structure is mathematically invariant with any static non-linguistic features. In short, phonology was implemented on physics, and the structure is purely acoustic and linguistic at the same time.

How to define the intelligibility quantitatively without any attention to the listeners? It is true that a student will communicate with many different non-native listeners and, in this meaning, the intelligibility may have to be defined based upon the listeners. However, it is true that the student is learning English of a single specific accent, i.e., British, American, Canadian, Australian English, or others. As is mentioned above, the pronunciations of individual students are acoustically and linguistically modeled as structures, which is similar to Halle's phoneme tree diagram [4]. It is also possible to extract the lexical structures from vocabulary of the individual Englishes. The pronunciation structure is determined by fixing a student and the lexical structure is determined by fixing an English. If compatibility between the two different levels of structures is measured, it will be another definition of the intelligibility. It is desired to measure the compatibility based upon some cognitive models because speaking is always intended to a human listener.

2. Physical implementation of phonology

2.1. Acoustic modeling of the non-linguistic information

Acoustic representation of speech with no dimensions to represent the static non-linguistic information in speech. How to derive it? To delete the non-linguistic information, it is modeled firstly, and then an algorithm for its deletion is implemented. In speech recognition, distortions caused by the non-linguistic events are often classified into three kinds; additive, multiplicative, and linear transformational. The additive distortion (noise) is ignored here because it is not inevitable. Students can turn off a TV set before learning English. The other two distortions are inevitable and their deletion has to be done by an algorithm.

Acoustic characteristics of microphones and rooms are typical examples of the multiplicative distortion. GMM speaker modeling indicates that a part of speaker individuality is also regarded as the multiplicative distortion. If a speech event is represented by cepstrum vector c , the multiplicative distortion is addition of b and the resulting cepstrum is shown as $c' = c + b$.

Vocal tract length difference is a typical example of the linear transformational distortion. The difference is often modeled as frequency warping of the log spectrum, where formant shifts are well approximated. According to [5], any monotonously continuous frequency warping of the log spectrum is mathematically converted into multiplication of matrix A in cepstrum domain. The resulting cepstrum is shown as $c' = Ac$.

Various distortion sources are found in every step of speech communication. But the total distortion of speech caused by the inevitable sources, A_i and b_i , is eventually modeled as $c' = Ac + b$, known as affine transformation.

2.2. From phonetics to phonology

In phonology, the non-linguistic information is mentally ignored in researchers' brain and speech sounds are represented as abstract entities named phonemes. Phonology clarifies a phonemic system or structure hidden in a set of the phonemes or in a sequence of the phonemes. Inspired by Saussure's structuralism, Jakobson, Halle, and others have discussed structure of a set of the phonemes embedded in a language with distinctive features and drew a tree diagram of the phonemes[4]. Classification of the phonemes is done so that a set of phonemes under every node of the tree comprise a natural class. In phonology, the structure is extracted in a top-down way based upon researchers' knowledge on the language. In this work, the structure is determined in a bottom-up way where not knowledge but distance measure between two elements is required. An n -point structure is represented uniquely by distance matrix among the n points. Viewing n elements as structure means that the elements are observed only relatively and the structure extraction can be regarded as a process of ignoring some information in the elements. If it is possible to embed all the sources of the non-linguistic information in the ignored information, the resulting structure will be the desired acoustic representation.

2.3. Implementation of phonology on physics

Phonology claims that the structure is universal with regard to all the kinds of non-linguistic information, which is mathematically translated that an n -point structure (distance matrix) is invariant with any affine transformation. This looks impossible, which can become possible by the following procedure.

Let phoneme x be represented as distribution $d_x(c)$ in a cepstrum space and distance between two elements (distributions) is calculated by Bhattacharyya distance (BD) measure.

$$BD(d_x(c), d_y(c)) = -\ln \int_{-\infty}^{\infty} \sqrt{d_x(c)d_y(c)} dc \quad (1)$$

This measure is derived based on information theory and can be interpreted as the amount of self-information of joint probability of the two independent distributions $d_x(c)$ and $d_y(c)$. If the two distributions follow Gaussians, the following is obtained.

$$BD(d_x(c), d_y(c)) = \frac{1}{8} \mu_{xy} \left(\frac{\Sigma_x + \Sigma_y}{2} \right)^{-1} \mu_{xy}^T + \frac{1}{2} \ln \frac{|\Sigma_u + \Sigma_v|/2}{|\Sigma_u|^{1/2} |\Sigma_v|^{1/2}} \quad (2)$$

μ_x and Σ_x are the average vector and the variance-covariance matrix of $d_x(c)$, respectively. μ_{xy} is $\mu_x - \mu_y$. Although affine transformation of $c' = Ac + b$ modifies $\mathcal{N}(\mu, \Sigma)$ into $\mathcal{N}(A\mu + b, A\Sigma A^T)$, BD between $d_x(c)$ and $d_y(c)$ is not changed.

$$BD(A\mu_x + b, A\Sigma_x A^T, A\mu_y + b, A\Sigma_y A^T) = BD(\mu_x, \Sigma_x, \mu_y, \Sigma_y) \quad (3)$$

These facts mean that BD between any two of the n distributions (phonemes) is not changed by any of an affine transformation and that the structure composed of the n phonemes is not changed. Multiplication of A and addition of b are geometrically interpreted as rotation and shift of the structure, respectively. For example, acoustic change of speech caused by

increase of vocal tract length, i.e., human growth, is mathematically regarded as very slow rotation of the structure which takes about 15 years. When $d_x(c)$ and $d_y(c)$ are modeled as Gaussian mixtures, the invariance is still valid because the structure of all the component Gaussians cannot be changed at all. Now, the desired acoustic representation is gracefully derived.

3. ERJ speech database

ERJ (English Read by Japanese) database[6] was used, which contains English read by 202 Japanese, Japanese English (JE), and 20 General American speakers (GA). The individual students have pronunciation scores rated by 5 American teachers of English. Table 1 shows the acoustic analysis conditions. Phoneme-to-phoneme distance is calculated as average distance over the three state-to-state BDs between two phoneme HMMs. Figure 1 shows a tree example of a Japanese student extracted from his HMMs and the well-known Japanese habits are clearly seen. Confusions of /t/ & /l/, /s/ & /th/, /z/ & /dh/, /f/ & /h/, /iy/ & /ih/, /v/ & /b/, etc are found. Mid and low vowels of English are located very close to each other because there is the only one mid and low vowel in Japanese. Schwa is close to the above vowels because Japanese often produce the mid and low Japanese vowel for schwa. Remarkably high performance of canceling the non-linguistic information was experimentally verified in [3, 7] and interested readers should refer to them.

4. Estimation of the intelligibility

In this section, compatibility between the pronunciation structure and the lexical structure is introduced based upon Cohort Model, one of the word perception models.

4.1. Cohort Model of word perception

The original Cohort Model characterizes a human process of perceiving an isolated word as a simple left-to-right process[8]. When the initial phoneme of a word input is perceived, a set of words starting with the phoneme are activated in brain. The

Table 1: Conditions for the acoustic analysis

sampling	16bit / 16kHz
window	25 ms length and 10 ms shift
parameters	FFT-based cepstrums and their derivatives
speakers	202 Japanese and 20 Americans
training data	60 sentences per speaker
HMMs	speaker-dependent, context-independent, and 1-mixture monophones with diagonal matrices
topology	5 states and 3 distributions per HMM
monophones	b,d,g,p,t,k,jh,ch,s,sh,z,zh,f,th,v,dh,m,n,ng,l,r,w,y,h,iy,ih,eh,ae,aa,ah,ao,u,h,uw,er,ax

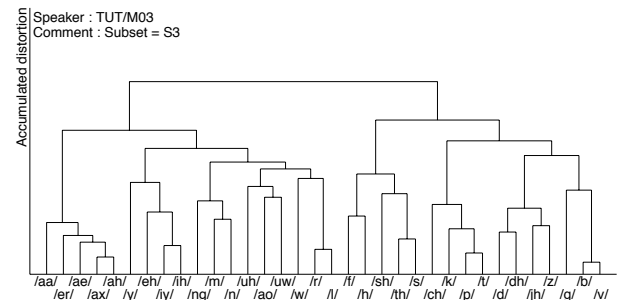


Figure 1: A structurally represented poor Japanese student

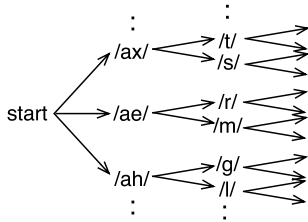


Figure 2: An example of the tree structured lexicon

number of the activated words is reduced by the subsequent input of phonemes and finally reaches one, which means the end of word perception. Cohort means a set of the activated words in brain. It is clear that Cohort Model assumes a tree structured lexicon in brain, which is shown in Figure 2. As is mentioned in Section 2.2, phonology clarifies a phonemic structure hidden in a set of the phonemes or in a sequence of the phonemes. The pronunciation structure discussed in the previous section corresponds to the former and is determined by fixing a student. The tree structured lexicon corresponds to the latter and is determined by fixing the target language. The intelligibility is defined as compatibility calculated between the two different levels of *phonological* structures based upon a cognition model. An algorithm for the calculation is shown below.

Cohort Model is often discussed with phonemes as its basic acoustic units. In this work, however, syllables are used as the basic units for cohort development. This is because an acoustic unit of speech production in English is said to be a syllable.

4.2. Estimation of the intelligibility as cohort size

Clearly seen in Figure 1, many phonemic confusions occur in Japanese English. This is natural because Japanese has only 25 phonemes and English has more than 40. If Japanese students use their own sounds only, 1-to- N mapping is inevitable. With the phonemic confusions, different words get acoustically closer and the acoustic lexical density is increased. In this work, larger lexical density is interpreted as less intelligibility. The compatibility between a student's pronunciation structure and the target language's lexical structure is defined as the cohort size calculated from the two structures. The smaller, the cohort size is, the higher, the compatibility and the intelligibility are.

The cohort activated only with the initial syllable input was focused upon. A 20K-sized lexicon in WSJ database was used as vocabulary and each entry of the lexicon has a unigram score. The phoneme sequence of each entry is obtained from PRONLEX dictionary. Each word (each phoneme sequence) was converted into a syllable sequence by tsylb software. Speech samples of some students in ERJ did not include a part of diphthongs and, in this case, HMMs for these phonemes could not be trained (See Table 1). Then, the words starting with a syllable including a diphthong as nucleus were ignored. The number of the remaining words was about 18K. It should be noted that the vocabulary includes different words whose baseforms are identical, such as walk, walked, and walking. Syllabification of the words showed that approximately 3,200 different kinds of syllables were found as word-initial syllables.

For each of the different word-initial syllables s_i , the number of the words starting with s_i or with a syllable acoustically close to s_i was calculated as $CS_0(s_i)$. Distance between two syllables was calculated by DP matching between two sequences of phoneme HMMs (syllables) and the calculation requires only the phoneme-to-phoneme distance matrix. The syllables acoustically close to s_i were defined as the syllables dis-

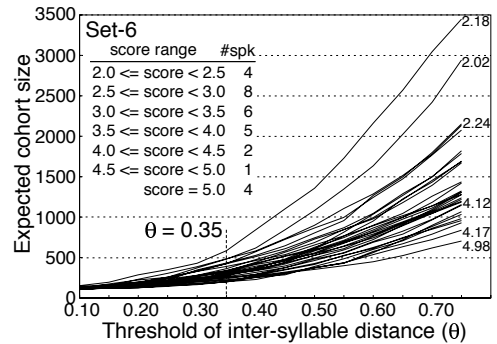


Figure 3: ECS as function of inter-syllable distance θ

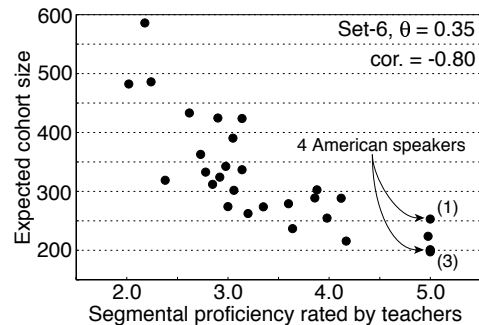


Figure 4: Proficiency rating without any acoustic matching

tant from s_i by less than threshold θ . Thus, $CS_0(s_i)$ was actually obtained as $CS_0(s_i, \theta)$, using which, size of the cohort activated by the initial syllable of word w_j was calculated as $CS_1(w_j, \theta) = CS_0(s^1(w_j), \theta)$, where $s^1(w_j)$ is the initial syllable of w_j . Finally, the expected cohort size $ECS(\theta)$ over the entire vocabulary was obtained by the following equation.

$$ECS(\theta) = \sum_j p(w_j) CS_1(w_j, \theta), \quad (4)$$

where $p(w_j)$ is a normalized uni-gram probability satisfying a condition of $\sum_j p(w_j) = 1.0$ over the words selected by deleting those starting with a syllable including a diphthong.

4.3. Results and discussions

Japanese students and GA speakers who read sentence set 6 were used in the experiment. The pronunciation structure somewhat depends upon the sentences read and set 6 was adopted because it covered a wide range of the proficiency with rather an even distribution (See Figure 3). The number of Japanese and Americans are 26 and 4, respectively. Proficiency scores of the Americans were assumed to be 5.0 (full score).

Figure 3 shows relations between the ECS and threshold θ for all the speakers, where the best three and the poorest three students are indicated by showing their pronunciation scores assigned by teachers. It is clearly indicated that words produced by the poorest students are very confusing and those by the best students are very distinct. This result shows good validity of the definition of the intelligibility adopted in this work.

Figure 4 shows correlation between the ECS and the pronunciation scores. Rather good correlation is found between the two. The ECS values are those at $\theta = 0.35$ in Figure 3 and the pronunciation scores were obtained by asking the teachers

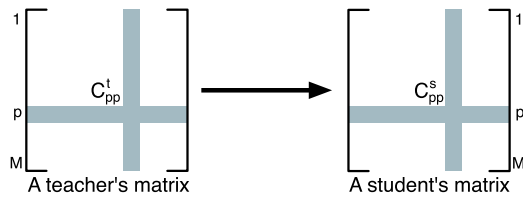


Figure 5: Replacement of a sub-structure

to rate the individual students with regard to the segmental aspect of the pronunciation. In the figure, the four Americans are explicitly indicated. Rather good correlation denotes high validity of the proposed method to estimate the intelligibility.

It should be noted that the proposed algorithm is implemented without any acoustic matching between a student and a teacher. The student's pronunciation is matched with the target language itself. The pronunciation structure can be said to be purely acoustic and linguistic at the same time. Then, the structure is matched with another level of linguistic structure, which is the lexical structure of the target language.

5. Effective and efficient instructions optimized for individual students

5.1. Exchange of sub-structures between two speakers

The pronunciation structure is extracted so that all the static non-linguistic features are discarded from speech. This characteristic enables an interesting operation, which is exchange of sub-structures between two speakers. If a sub-structure in a student is replaced with the corresponding one in a teacher, the student will have a better pronunciation structure (See Figure 5). This operation is meaningless if the other acoustic representation of speech, spectrogram, is used. If a portion of the spectrogram of a speaker is replaced with the corresponding spectrogram of another, all the acoustic features are changed, such as proficiency, age, gender, size, shape, microphone, room, line, and so on. The resulting spectrogram has to be just a mess.

5.2. Instructions optimized for individual students

Replacement of which sub-structure minimizes the cohort size? The answer to this question will provide the pedagogical instruction optimized for the student. In the current work, a sub-structure of phoneme p is defined as the following set of elements in the distance matrix; $\{c_{pj}\}$ and $\{c_{jp}\}$. Here, c_{ij} is an element of the matrix. If replacement of the sub-structure of phoneme p_0 minimizes the cohort size, it means that the student should correct the articulation of phoneme p_0 among others.

Using a female speaker, RYU/F06 (pronunciation score is 2.02), the cohort size reduction is done by replacing her sub-structures with a teacher's ones. Figure 6 shows results of the cohort reduction by a single replacement. Her original cohort is more than double of the native cohort and replacement of a sub-structure of schwa is shown to be the most effective and efficient for her to improve the intelligibility of the pronunciation. What's next if schwa is corrected? The most effective replacement after schwa's correction can be discussed in the same manner. Figure 7 shows the order of English phonemes for her to correct and size of the effect accumulated by the sequential corrections. It is shown in Figure 7 that the phonemes with higher priority for their corrections are those which are known to be pronounced by Japanese to be acoustically similar to other phonemes. This result shows good validity of the proposed method for automatic generation of instructions.

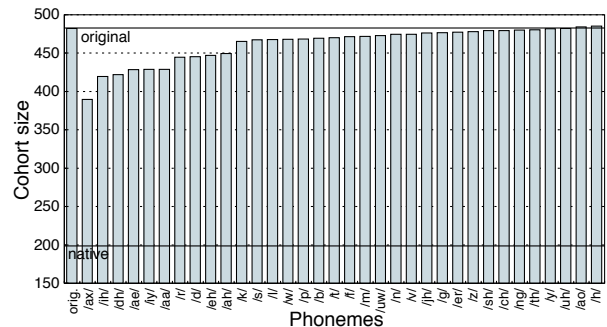


Figure 6: Cohort size reduction by a single replacement

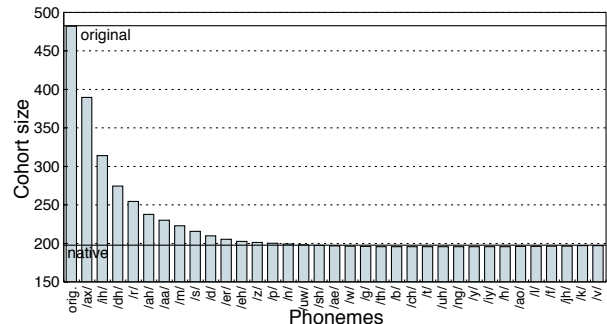


Figure 7: Cohort size reduction by sequential replacements

6. Conclusions

This paper proposed a novel method to estimate compatibility between a student's pronunciation structure and the target language's lexical structure, which is regarded as the intelligibility of the pronunciation. The proposed algorithm does not require any acoustic matching between a student and a teacher, which means that the algorithm cannot face "mismatch problems" at all. This paper also showed that it is possible to determine the order of phonemes for individual students to correct. This determination is done by sequential replacements of sub-structures and this operation is possible only with the acoustic and linguistic representation of speech which the author proposed previously. As future work, the author is planning to verify the effectiveness of the proposed method in actual classrooms with not only university students but also young children.

7. References

- [1] J. Bernstein, "Objective measurement of intelligibility," Proc. ICPhS, pp.1581–1584 (2003)
- [2] N. Minematsu *et al.*, "CART-based factor analysis of intelligibility reduction in Japanese English," Proc. EUROSPEECH, pp.2069–2072 (2003)
- [3] N. Minematsu, "Yet another acoustic representation of speech sounds," Proc. ICASSP, pp.585–588 (2004)
- [4] M. Halle, "The sound patterns of Russian," The Hague: Mouton (1959)
- [5] M. Pitz *et al.*, "Vocal tract normalization as linear transformation of MFCC," Proc. EUROSPEECH, pp.1445–1448 (2003)
- [6] N. Minematsu *et al.*, "Development of English speech database read by Japanese to support CALL research," Proc. ICA, pp.557–560 (2004)
- [7] N. Minematsu, "Pronunciation assessment based upon the phonological distortions observed in language learners' utterances," Proc. ICSLP (2004, submitted)
- [8] W. D. Marslen-Wilson *et al.*, "The temporal structure of spoken language understanding," Cognition, vol.8, pp.1–71 (1980)