

BRITISH ACOUSTICAL SOCIETY

"SPRING MEETING" at Chelsea College, London, S.W.3. on
Wednesday 25th April / Friday 27th April, 1973.

SPEECH AND HEARING: Session 'C' Speech Properties and Recognition.

Paper No:

73SHC2

On the Adaptation in Speech Recognition

B.M.LOBANOV

Radiotechnical Institute of Minsk, USSR

(Invited paper)

It is well known that the speech signal is extremely variable, depending on the characteristics of the acoustical tract, vocal peculiarities of the speaker and his physiological and psychological state. Due to the slow changing nature of these interfering factors an adaptation procedure of speech recognition is made possible.

In the present paper an attempt of working out an adaptation procedure model with the object of automatic recognition on the basis of the available data on adaptation and the authors previous investigations is made.

Let $\{\Phi_1, \dots, \Phi_j, \dots, \Phi_N\}$ be a set of speech patterns intended for recognition; $\{S_1, \dots, S_K, \dots, S_M\}$ a set of the speech signal parameters and $\{\Pi_1, \dots, \Pi_i, \dots, \Pi_P\}$ a set of interfering factors conditioned by the influence of the individual characteristics of the electroacoustical tract and the speaker. Automatic recognition of a given set of speech patterns is only possible if the chosen speech signal parameters meets definite requirements. The minimal requirement is that the distributions of the speech signal parameters for each of the speech patterns should not intersect (above a preliminary fixed level) at any fixed interfering effects.

A set of speech signal parameters may be considered optimal if they are fully invariable in respect of any interfering effects (labeled a priori adapted parameters). It means that the speech pattern distribution centers vary so insignificantly that it is possible to

apply the same decision rule at different values of interfering factors. In this case adaptation procedure is not necessary. Unfortunately, an invariant system of speech parameters (or features) can only be found but for a limited set of speech units.

Radically different is a system of speech signal parameters with the distribution centres varying considerably in different a priori unknown directions at variations of characteristics of interfering factors (labeled fully unadapted parameters). Adaptation is quite necessary then and a complete readjustment of all elements of the recognition system is required with the help of displaying all the speech patterns intended for identification.

Between these two is a set of parameters having the same exact or statistical linear regularities in the distribution centres displacement for all the speech patterns (labeled linear unadapted parameters). If an interfering effect produces a linear transformation of the space of parameters reduced to expansion (or compression) and shifting along the coordinate axes, then we have equation $S_{kj}^i = a_k^i S_{kj}^0 + b_k^i$, (1) where S_{kj}^i is the value of the k^{th} parameter for the j^{th} speech pattern at the i^{th} state of the interfering factors; S_{kj}^0 is the value of the k^{th} parameter for j^{th} the speech pattern if interfering factors are not present; and a_k^i and b_k^i are expansion and shifting coefficients. In this case adaptation process may be accomplished after displaying of only two realizations of preliminary known different patterns Φ_α and Φ_β . The unknown a_k^i and b_k^i are determined from the system of two linear equations

$$\begin{cases} S_{k\alpha}^i = a_k^i S_{k\alpha}^0 + b_k^i \\ S_{k\beta}^i = a_k^i S_{k\beta}^0 + b_k^i \end{cases}$$

and then adapted values of the parameters are found according to the formula

$$S_k^A = \frac{S_k^i - b_k^i}{a_k^i} = \frac{S_k^i (S_{k\beta}^0 - S_{k\alpha}^0) - (S_{k\alpha}^i S_{k\beta}^0 - S_{k\beta}^i S_{k\alpha}^0)}{S_{k\beta}^i - S_{k\alpha}^i}. \quad (2)$$

The unknown coefficients a_k^i and b_k^i can also be found directly during the process of recognition that is by means of self-training. For this purpose it is sufficient to define some extremal [1] or mean-statistical values [2] of parameters. Really, substituting in formula (2)

$S_{k \max}^i$ and $S_{k \min}^i$ for $S_{k\beta}^i$ and $S_{k\alpha}^i$, assuming that $S_{k \max}^0 = 1$ and $S_{k \min}^0 = 0$ we come to the formula $S_k^A = \frac{S_k^i - S_{k \min}^i}{S_{k \max}^i - S_{k \min}^i}$, suggested by Gerstman [1]. Considering, further, in formula (1), S_{kj}^i to be a random value with $M=0$ and $\sigma=1$ and finding for S_{kj}^i the mean and root mean square values we can write

$$M(S_{kj}^i) = M(a_k^i S_{kj}^i + b_k^i) = b_k^i$$

$$\sigma(S_{kj}^i) = \sigma(a_k^i S_{kj}^i + b_k^i) = a_k^i$$

that is S_k^A according to formula (2) can be expressed by the formula

$$S_k^A = \frac{S_k^i - M(S_{kj}^i)}{\sigma(S_{kj}^i)},$$

suggested by the author [2].

Information of the expansion and shifting coefficients can also be obtained through using ^{their} statistical correlations with the other parameters that are not directly involved in the process of recognition. It has been shown [3] that current values of highest formants and fundamental frequencies can be used in adaptation process.

We have noted, besides, that on all occasions a_k^i and b_k^i have to be determined only for those speech patterns for which they are identical. In general the set of speech patterns may be subdivided into a number of groups [4], for which a_k^i and b_k^i are determined separately.

Different types of adaptation should be applied in speech recognition depending on the parameter system employed. Let us consider, for instance, two parameter systems: a) signal-parameters of the parallel spectral analyzer and b) formant frequencies of the articulatory tract. The first parameter system may be considered "a priori adapted" at buzz and hiss sounds identification. But for the majority of the phonetic elements the system of signal parameters is either "linear unadapted" (for interfering factors due to ^{changing of the} frequency characteristics of the electroacoustical tract and spectral characteristics of the excitation sources) or "fully unadapted" (for interfering factors due to the ^{average} changing of the sizes of the speaker's articulatory tract). The second system is better since it is "a priori adapted" for the first type of interfering factors and only "linear unadapted" for the second type of interfering factors.

In temporal aspect we could offer the following organisation of the adaptation procedure. At the first stage (a priori adaptation) the average characteristics of the most probable type of the speaker and the electroacoustical tract must be taken. At the second stage (momentary adaptation) the average characteristics are corrected by taken into account the information about the current values of some speech parameters (the fundamental frequency, the frequency of the third formant, etc.). At the final stage (integral adaptation) they are finally corrected by calculating the statistical characteristics of the speech parameters on the chosen time intervals.

The author is indebted to L. Karnevskaya and L. Leladze for the help in preparation of the English version of the present paper.

REFERENCES

1. L.J. Gerstman. Classification on Self-Normalized Vowels. IEEE Trans. on Audio and Electroacoust., v. AU-16, No1, 1968.
2. B.M. Lobanov. Classification of Russian Vowels Spoken by Different Speakers. J.A.S.A., v.49, No2 (p.2), 1971.
3. H. Fyjisaki, T. Kawashima. The Roles of Pitch and Higher Formants in the Perception of Vowels. IEEE Trans. on Audio and Electroacoust., v. AU-16, No1, 1968.
4. B.M. Lobanov. On the Classification of Russian Fricatives in the CV-Syllables for Different Speakers. J.A.S.A., v.49, No4 (p.2), 1971.