



Analysis of Autocorrelation-based Parameters for Creaky Voice Detection

Carlos Toshinori Ishi

JST/CREST,
ATR/HIS Labs. Kyoto, Japan
carlos@atr.co.jp

Abstract

Creaky voice carries important linguistic and paralinguistic information. Parameters based on autocorrelation of the glottal excitation waveform are proposed for automatic detection of creaky voice in spontaneous speech. Analysis results show the ratio of the first two peaks of the autocorrelation function as a primary parameter to detect creaky voice.

1. Introduction

“Creaky voice” has many other terminologies, such as “creak”, “vocal fry”, “glottal fry”, “laryngealization”, “glottalization”, and “pulse register phonation”, used in several research areas (like Linguistics, Physiology, and Phonetics) [1,2].

Creaky voice is defined as “... a train of discrete laryngeal excitations, or “pulses”, of extremely low frequency (7 to about 78 Hz), with almost complete damping of the vocal tract between excitations.” (Hollien 66 cited in [1]). “The auditory effect is of a rapid series of taps, like a stick being run along a railing.” (Catford 64 cited in [2]).

Creaky phonation carries many linguistic and paralinguistic information, depending on the language. For example, contrast between creaky and modal voicing among vowels and nasals is particularly common in some American Indian languages [3]. In [4,5], relationship between different phonation types and paralinguistic information like emotions and attitudes are reported for English. Strong correlations were reported between creaky voice and perception of relaxed/stressed, sad/happy, and bored/interested. In Japanese, expressive pressed voice that is frequently realized by creaky phonation also carries important paralinguistic information such as attitudes, emotional states and emphasis [6].

Further, in creaky segments, periodicity is disturbed and the pitch extraction becomes difficult, affecting the subsequent prosodic analysis, like intonation. Tendency of creaky segments for specific tone types is reported in [7] for phrase finals in Japanese.

In the JST/CREST ESP Project [8], one of the goals is an expressive speech synthesizer based on unit selection, using a large database of spontaneous speech. For this purpose, labels of voice qualities (phonation types) become as important as prosodic labels. With the goal of doing automatic labeling of voice quality on a large speech database, in the present research, we focus on the automatic detection of creaky phonation.

2. Acoustic features of creaky voice

A lot of research has been conducted on the acoustical analysis of creaky voice. Among them, we can cite disturbance of periodicity in the time domain (jitter and shimmer) [11], feature extraction in the power spectrum [3,10], and parameterization of the glottal excitation waveform obtained from vocal tract inverse filtering of the speech signal [4,9].

Jitter (perturbation in the fundamental frequency) and shimmer (perturbation in the amplitude) are two measures of perturbation in the periodicity in the time axis. There are many works [11] showing their correlation with perceptual roughness. However, direct correlation with creaky phonation is not reported, since creak can also be periodic.

[3] describes that creaky phonation has the following features compared to modal phonation: non-periodic glottal pulses, lower power, lower spectral slope, low F0. Among them, the spectral slope is reported to be the most important parameter to discriminate between different phonation types. In [3], the spectral slope is estimated based on harmonic components of the power spectrum. [10] also estimates the spectral slope based on harmonics of the spectrum, but considering the effects of the formants. However, this kind of method using harmonic components could not be appropriate for non-periodic signals.

Another approach for discriminating phonation types is based on speech production. The basic idea is removing the effects of the vocal tract resonances from the speech signal by inverse filtering techniques, to obtain the glottal excitation waveform. In the research field of speech synthesis based on speech production models, the glottal excitation waveform is parameterized according to the shape of each glottal pulse. [4,9] reports successful synthesis of different voice qualities, including creaky voice, by controlling the parameters of the LF model.

However, automatic detection of creaky voice is not as widely reported. Perhaps because the glottal excitation is irregular and automatic detection becomes difficult.

3. Autocorrelation-based parameters

In the present research, in order to avoid the detection of excitation pulses in the temporal domain, we propose a parameterization of phonation type features based on the autocorrelation of the glottal excitation waveform.

3.1. Estimation of the glottal excitation waveform

The estimation of the glottal excitation waveform is based on

the method proposed in [12]. First, the speech signal is high-pass filtered at 60 Hz in order to prevent the glottal waveform from gradually rising (or falling). Then, the glottal contribution to the speech spectrum is preliminarily estimated by applying LPC-analysis of order 1. We refer to the estimated coefficient as the adaptive pre-emphasis (*APE*) coefficient. Next, the speech signal is pre-emphasized by using *APE*, and LPC-analysis of order 18 (with sampling frequency of 16000 Hz) is applied on the pre-emphasized signal. The obtained LPC coefficients are used for inverse filtering of the high-pass filtered speech signal. The residual signal is treated as the glottal excitation waveform hereinafter.

3.2. Normalized Autocorrelation Function (*NACF*)

Before estimating the autocorrelation function (*ACF*), the glottal excitation waveform is low-pass filtered at 2 kHz, in order to make the *ACF* peak detection easier.

An important point to be taken into account is the window size for *ACF* estimation. Since creaky voice usually appears in low fundamental frequencies, the window size should be long enough to cover at least two excitation pulses. On the other hand, a too long window size is not appropriate for segments with high and changing pitch. Therefore, we decided to use an analysis window with variable length.

A two-step *ACF* estimation is used to adjust the window length adaptively. First, *ACF* is estimated in an 80 ms window. And then, the time lag of the maximum peak is extracted and multiplied by four, to be used as the new window size. Here, the new window size was clipped to lie in the interval between 16 ms and 80 ms.

The obtained *ACF* is normalized according to the following expression:

$$NAC(L) = \frac{N}{N-L} \frac{R_{xx}(L)}{R_{xx}(0)}, \quad (1)$$

where N is the number of samples of the frame window, L is number of samples of the autocorrelation lag, and R_{xx} is the autocorrelation function. Figure 1 shows examples of glottal excitation waveforms and normalized autocorrelation functions

obtained using the methods described above.

For modal phonation (a), a clear periodicity can be observed; the *NACF* peaks are close to 1 value, and there are no small peaks between the time lag 0 and the first big peak. (b) and (c) show examples of creaky voice with big-small-big-small and short-long-short-long sequences (jitter/shimmer) of the glottal pulses. (b) shows a smaller peak between the time lag 0 and the maximum peak. The magnitude of this small peak becomes lower and the width of this peak becomes larger, as the jitter/shimmer becomes stronger, such that it divides in two, as shown in (c). In the modal phonation example, it can also be observed that the first two peaks (closer to time lag 0) have values close to 1. (d) shows an example of (non double-periodic) creak, where only one big *NACF* peak can be observed. However, a narrow width is observed for this peak, because of the impulse-like shape of the glottal excitation for creak phonation.

3.3. *NACF*-based parameters

Based on visual inspection of the *NACF* of the glottal excitation waveforms of modal and creaky phonations as described in the previous section, we decided to use the first two peaks (called $P1$ and $P2$, from the time lag 0) in the *NACF*, to characterize different phonation types. A threshold of 0.2 is used to detect peaks in *NACF*. The following parameters are proposed based on these two peaks ($P1$, $P2$).

- **Peak magnitude (*NAC* value) ratio:**

$$NACR = 1000 * NAC(P2) / NAC(P1) \quad (2)$$
- **Peak position (time lag) ratio:**

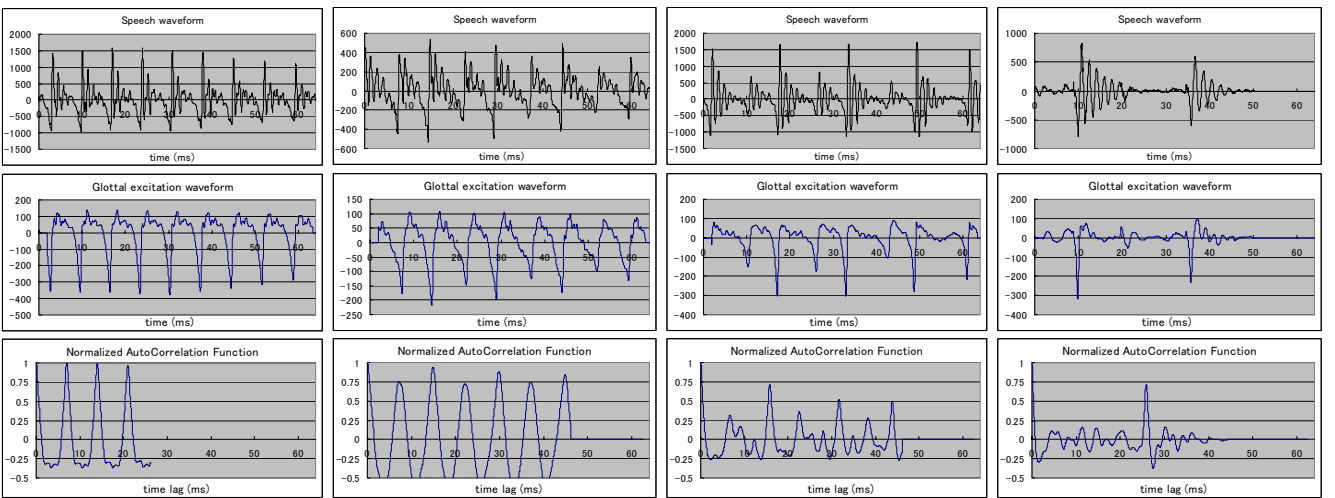
$$TLR = 2000 * TL(P2) / TL(P1) \quad (3)$$
- **Peak width ratio:**

$$WR = 1000 * W(P2) / W(P1) \quad (4)$$
- **Maximum peak magnitude:**

$$NACmax = 1000 * NAC(Pmax) \quad (5)$$
- **Maximum peak position:**

$$TLmax = TL(Pmax) \quad (6)$$
- **Maximum peak width:**

$$Wmax = W(Pmax) \quad (7)$$



(a) modal phonation (b) low jitter creaky phonation (c) high jitter creaky phonation (d) low F0 creaky phonation

Figure 1: Speech waveform, glottal excitation waveform and *NACF* for Modal and Creaky phonation

A scaling factor of 1000 is multiplied in the $NACR$, TLR , WR , and $NACmax$ parameters for data storage convenience, allowing their storage in integer format. Tables 1 and 2 show the expected behavior of the proposed parameters for each phonation type. For single-periodicity, all the ratios ($NACR$, TLR , WR) and $NACmax$ are expected to have values close to 1000. For double-periodicity, $NACR > 1000$; $NACmax < 1000$; if jitter is strong, $TLR \neq 1000$; and if jitter or shimmer is strong, $WR < 1000$.

For low F0 creaky phonation with non-double periodicity, i.e., large interval between excitation pulses (big $TLmax$), there are cases where only one peak can be detected. Therefore, the ratio-based parameters cannot be used to represent these signals. However, a small value of $Wmax$ is expected in these cases, since creaky phonation has narrow (impulse-like) excitation pulses (Table 2).

Table 1: Expected behavior of the parameters in modal and double periodic signals.

	$NACR$	TLR	WR	$NACmax$
(Single) Periodicity Modal	$\cong 1000$	$\cong 1000$	$\cong 1000$	$\cong 1000$
Double Periodicity Creaky/Rough	> 1000	$\neq 1000$	< 1000	< 1000

Table 2: Expected behavior of the parameters in low fundamental frequency creaky phonation.

	$TLmax$	$Wmax$
Low Frequency Creaky	Big	Small

4. Analysis and evaluation of the proposed parameters

As speech data for evaluation, we used the same dataset analyzed in [7], containing 404 phrase final syllables

segmented from natural spontaneous speech of a female adult speaker. Each syllable was labeled in terms of Creaky(C), Aspirated(A) or Modal(M), looking at the waveform and hearing the segments. The parameters proposed in Section 3 were estimated in all frames (5619) of the annotated speech intervals.

As a preliminary evaluation, a decision tree was constructed for each of the categories {C,A,M}, using the R Package[12]. The tree resulted in 91.5% of correct identification. Specifically for Creaky category, deletion error was 13.7%, while substitution error was 7.9%. However, only the parameter set { $NACmax$, $NACR$, TLR } was used in the constructed decision tree. A detailed analysis of each parameter was then conducted to verify their behavior in each category.

Figure 2 shows the distributions of each parameter separated for each voice quality category. The data of all panels are arranged according to increasing order of $NACR$. It can be observed from the panels that the data of $NACR$, TLR , WR and $NACmax$ for Modal category are concentrated around the value 1000, as expected in Table 1. In order to better visualize the distributions of each category, 3155 frames of the Modal category in the interval $935 < NACR < 975$ were removed from the panels, since the distributions of all parameters showed regularity relative to the adjacent portions.

Dashed lines were placed separating the regions where $NACR < 1000$ and $NACR > 1000$, for the categories C (C2: 268 frames, and C3: 630 frames) and M (M1:961+3155 frames and M2: 86 frames), resulting in a clear separation of data in both WR and $TLmax$ parameters. The features of C2 are similar to those of M1, while the features of M2 are similar to those of C3. Part of C2 and M2 data are frames close to the boundaries between creaky and modal regions. Also part of M2 data are frames close to the boundaries between modal and silence, where laryngealizations frequently occur. $TLmax$ is the autocorrelation time lag of the maximum peak, and this value is also a base value for F0 estimation. It is clear that the sets

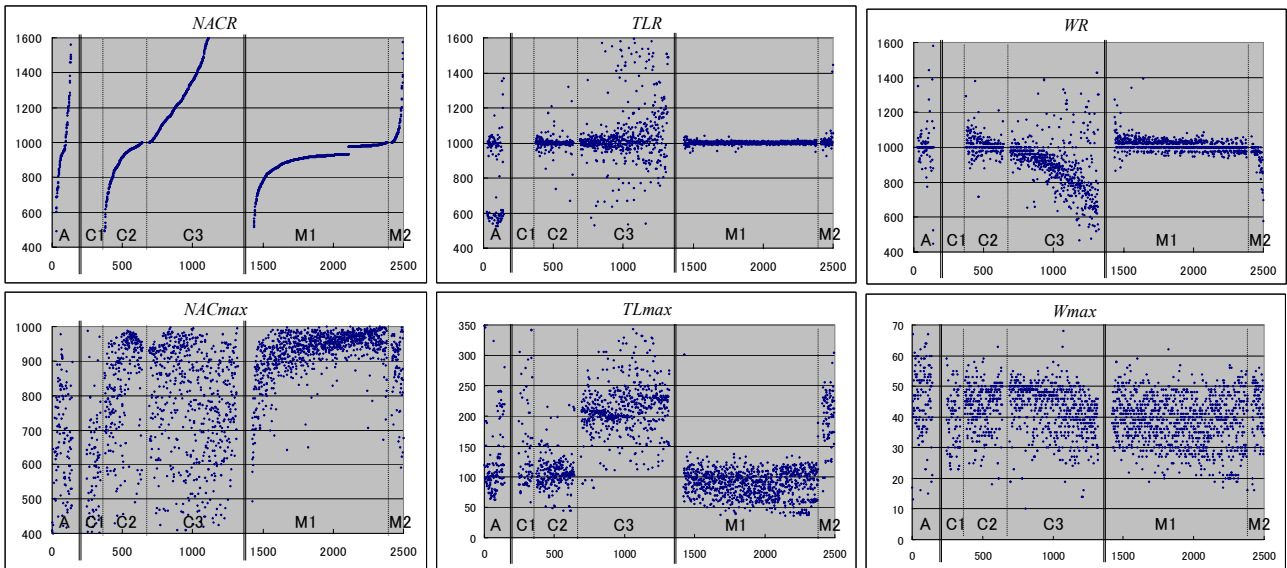


Figure 2: Distributions of autocorrelation-based parameters for each voice quality category (Aspirated, Creaky, and Modal). Abscissas are number of frames. Ordinates of $NACR$, TLR , WR and $NACmax$ have values scaled by 1000.

{C2,M1} and {C3,M2} have distinct distributions in the $TLmax$ panel. This is probably because of double-periodicity in {C3,M2}. This result can also be used for double-pitch error correction, since double-periodicity consistently occurs when $NACR > 1000$.

A strong inverse correlation can be observed between $NACR$ and WR . This could be a reason why WR was excluded by the decision tree algorithm.

The subgroup C1 (84 frames), in the category C, represents the frames where only one peak was detected; this means that no data can be plotted in the ratio-based parameters ($NACR$, TLR and WR). $NACmax$ shows lower values for C1 compared to M.

Unfortunately, no clear separation can be noted in the distributions of $Wmax$ parameter. It is expected that more significant differences will appear for male speaker voices, especially between creaky voice and lax voice with low fundamental frequencies. Analyses are also being conducted for male speakers.

No clear discrimination can be observed between Creaky and Aspirated categories except for the TLR parameter. However, only 33% (47/143) of the Aspirated frames can be separated by using a threshold of 650 for TLR . As a solution to improve the discrimination between these categories, we propose that the coefficient of the adaptive pre-emphasis (APE , described in Section 3.1) can be used as a distinctive parameter, since aspirated speech intervals tend to be stronger in the high frequency bands, resulting in a lower magnitude for the pre-emphasis coefficient. Figure 3 shows the distribution of APE for each category.

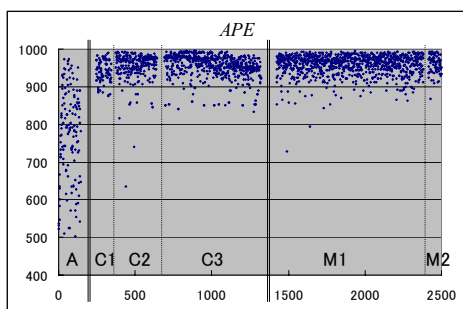


Figure 3: Distributions of Adaptive Pre-emphasis coefficient (APE) parameter.

72% (102/143) of the Aspirated frames are correctly separated by using an APE threshold of 800. 84% (121/143) can be reached if the APE threshold is set to 900. Other parameters should be evaluated in the discrimination between Creaky and Aspirated phonations.

Removing the vocal tract effect (by inverse filtering) to obtain glottal excitation is a convenient approach to analyze phonation types from a physiological viewpoint. However, from a perceptual viewpoint, it is hard to say that the human ear realizes LPC inverse filtering to separate between the glottal source and vocal tract components. From this perspective, an auditory model could be used to evaluate phonation type identification without applying inverse filtering on the input speech signal.

5. Conclusion

Parameters based on the normalized autocorrelation function of glottal excitation waveform were investigated with the aim of automatically detecting creaky voice segments. Preliminary evaluation of the proposed parameters showed good performance in the automatic detection of creaky voice. Among the parameters, the ratio between the first two autocorrelation peaks ($NACR$) was found to be the primary parameter to discriminate between modal and creaky phonation. Once creaky segments are detected, the next step is to verify if these segments are really perceived as rough. Another topic to be investigated is how pitch is perceived in the creaky intervals.

6. Acknowledgements

I would like to thank Nick Campbell and the whole JST/CREST group for supporting the present research. Special thanks for Parham Mokhtari (ATR) and Ken-ichi Sakakibara (NTT) for advice and motivating discussions.

7. References

- [1] Gerratt, B. R., Kreiman, J., 2001. Toward a taxonomy of nonmodal phonation. *J. of Phonetics* 29, 365-381.
- [2] Laver, J., 1980. Phonatory settings. In *The phonetic description of voice quality*. Cambridge University Press, Ch. 3, 93-135.
- [3] Gordon, M., Ladefoged, P., 2001. Phonation types: a cross-linguistic overview. *J. of Phonetics* 29, 383-406.
- [4] Gobl, C., Ni Chasaide, A., 2003. The role of voice quality in communicating emotion, mood and attitude. *Speech Communication* 40, 189-212.
- [5] Klasmeyer, G., Sendlmeier, W. F., 2000. Voice and Emotional States. In *Voice Quality Measurement*, Singular Thomson Learning. Ch. 15, 339-358.
- [6] Sadanobu, T. 2003. Expressive Speech and Grammar: with special reference to pressed voice in Japanese. *JST/CREST Int. Workshop on Expressive Sp. Proc.*, 55-60.
- [7] Ishi, C.T., Mokhtari, P., Campbell, N. 2003. Perceptually-related Acoustic-Prosodic Features of Phrase Finals in Spontaneous Speech. *Eurospeech 2003*, 405-408.
- [8] The JST/CREST Expressive Speech Processing project, introductory web pages at: www.isd.atr.co.jp/esp
- [9] Childers, D.G. 1995. Modeling the glottal volume-velocity waveform for three voice types. *J. Acoust. Soc. Am.* 97 (1), 505-519.
- [10] Hanson, H. M., Stevens, K., Kuo, H. J., Chen, M., Slifka, J., 2001. Towards models of phonation. *J. of Phonetics* 29, 451-480.
- [11] Buder, E.H. 2000. Acoustic Analysis of Voice Quality: A Tabulation of Algorithms 1902-1990. In *Voice Quality Measurement*, Sing. Thomson Learning, Ch. 9, 119-244.
- [12] Alku, P., Vilkmann, E., Laine, U. 1990. A comparison of EGG and a new automatic inverse filtering method in phonation change from breathy to normal. *ICSLP 1990* (1), 197-200.