

The effects of anger on automated long-term-spectra based speaker-identification

Ortega-Rodríguez, M.¹ ; Solís-Sánchez, H.¹ ; Valverde-Méndez, D.^{1,2} ; Venegas-Li, A.^{1,3} 

¹ School of Physics and Geophysical Research Center, University of Costa Rica, San José, Costa Rica, {manuel.ortega, hugo.solis}@ucr.ac.cr

² Department of Physics, Princeton University, Princeton, United States, dsmendez@princeton.edu

³ Physics Department, University of California at Davis, Davis, United States, avenegasli@ucdavis.edu

Abstract

Forensic speaker identification has traditionally considered approaches based on long-term (a few tens of seconds) spectra analysis as especially robust. This is because they work well for short recordings, are not sensitive to changes in the intensity of the sample, and continue to function in the presence of noise and limited passband. Because of this, the long-term spectra approach is one of the preferred tools for forensic speaker identification, in addition to formant analysis, speed of speech, and determination of the fundamental frequency. However, we find that anger induces a significant distortion of the acoustic signal for long-term spectra analysis purposes. Even moderate anger offsets speaker identification results by 33% in the direction of a different speaker altogether (in the space of sample correlations). Therefore, caution should be exercised when applying this tool.

Keywords: automated speaker identification, long-term spectra, forensic acoustics, emotional distortions, anger.

PACS: 43.72.Uv, 43.72.Ar, 43.72.Fx.

El efecto del enojo en los procesos automatizados de identificación forense de personas locutoras basados en espectros del habla a largo plazo

Resumen

La identificación forense de locutores/locutoras ha considerado tradicionalmente acercamientos al problema basados en el análisis de espectros a largo plazo (varias decenas de segundos de duración). Estos acercamientos han demostrado ser especialmente robustos, en el sentido que siguen funcionando bien incluso si las grabaciones son cortas; además, el método no es sensible a cambios en la intensidad sonora de la muestra, y sigue funcionando bien en la presencia de ruido y de ancho de banda limitado. Por todo esto, constituye una de las técnicas preferidas para la identificación forense, junto con el análisis de formantes, la velocidad del habla y la determinación de la frecuencia fundamental. Se halla, sin embargo, que el estado de enojo produce una distorsión importante en la señal acústica para efectos del análisis de espectros del habla a largo plazo. Incluso si el nivel de enojo es solamente moderado, hay un desvío de los resultados cuantitativos de la identificación forense de personas locutoras que representa el 33 % de la distancia (en el espacio de correlación entre muestras) hacia una persona locutora totalmente distinta. Por tanto, se concluye que es importante tener cautela en el momento de aplicar este método.

Palabras clave: identificación forense de locutor y locutora, espectros a largo plazo, acústica forense, distorsiones emocionales, enojo.

1. INTRODUCTION

The purpose of this article is to study how to quantitatively determine the effect of distortions caused by emotional states (particularly anger) on the analysis of long-term speech spectra (also known as LTS, *long-term spectra*) for Forensic Speaker Identification purposes (FSI). (Examples of articles by experts in the field of FSI are given by Hollien [1] and Hollien [2].) This research is conducted using a careful and reproducible methodology, which is explained below. The aim of the FSI process is to identify a speaking person through the analysis of their voice, usually under less-than-ideal conditions. One of the fundamental challenges FSI must face is determining whether the within-speaker variability is less than the between-speaker variability (which is clearly desirable) and how this relationship holds under different conditions (Hollien [3]). Among the most common conditions are technological distortions due to the equipment used for recording, as well as environmental distortions caused by noise or harsh background sounds.

In particular, speakers themselves can be a source of distortion, as there can be a variety of feelings such as fear, anger, and anxiety (a likely situation, for example, in FSI when the speaker might be committing a crime). These emotions trigger a modification in speech production that manifests itself as a change in the values of signal parameters (such as frequencies and speech rate) (Williams and Stevens [4]; Banse and Scherer [5]; Johnstone [6]). Voice production consists of air pulses caused by the vibration of the vocal cords (which are then modified by the supralaryngeal vocal tract), so the dominant factors of vocalization are breathing patterns and the varying tension of the muscles involved in the process. Since these factors are highly correlated with emotions, it is therefore very likely that such changes in the factors will be detectable in the acoustic wave (Scherer [7]).

However, this topic has not been extensively studied for any language, largely because there are ethical and methodological restrictions that

complicate the controlled production of strong emotions. The consensus (Johnstone [6]) has been that controlled laboratory conditions are only feasible for low-intensity emotional states, for which it is rather difficult to notice a change in the effectiveness of FSI. Another difficulty lies in how to induce the desired emotion. Martin [8] provides an overview of some possible techniques used to generate specific emotions, such as music and emotive images, as well as self-generation techniques such as the use of imagination and memories. These methods are classified according to the way emotions are produced.

For the purposes of this research, the method of self-induced autobiographical memories was selected. In this technique, participants (all men, as explained in more detail below) are asked to recall emotive events in order to generate the desired emotion. Although this is not the only method applicable to the problem at hand, it was chosen for its simplicity and because it allowed participants to have privacy while making the recordings. Furthermore, the quality of a blind experiment was ensured by having the participants themselves determine their level of anger (on a numerical scale), as described later.

2. BACKGROUND

This section will discuss the logic behind the process of forensic speaker identification through the use of long-term spectral analysis (LTS), a process that will henceforth be referred to as FSILTS. Although there are many markers (also known as “vectors”) that help distinguish recordings of different speakers since the pioneering work of Hollien and Majewski [9], one of the most common methods employed in FSI is LTS analysis (Kinnunen *et al.* [10]; Ortega-Rodríguez *et al.* [11]). LTS analysis quantitatively reveals the temporal average of the timbre of the voice, which is the acoustic property that allows a listener to distinguish between a clarinet and a violin playing the same musical note (frequency) at the same sound intensity, for example. The distribution is obtained by computing the Fourier Transform of the signal over a long period, such as 30 seconds. The idea is that

the speaker has enough time to move through the entire sound phase space, meaning that they have the opportunity to pronounce several times all (or almost all) the sounds of Spanish.

This vector has been extensively studied in terms of the efficiency with which it is able to identify the speaker. As a result, it has been found to be one of the most reliable, mainly because it continues to function even in the presence of noise and limited bandwidth (Hollien [3]).

One of the challenges when using this vector is how to define the correlation between two spectra. There are several ways to approach this problem. Among the most common are assigning a specific number to each LTS data set (according to some algorithm), or even visual inspection of the graphs. However, for the purposes of this article, a more sophisticated method is needed. Two correlation coefficients were considered for these purposes: the Standard Deviation of the Differences Distribution (SDDD) (Harmegnies [12]) and the Bravais-Pearson cross-correlation coefficient, R (Stanton [13]). Several exploratory experiments carried out by our group without the anger component (Ortega-Rodríguez *et al.* [11]) showed that the Bravais-Pearson correlation coefficient gives the best results for this line of research and was therefore selected. (The criterion used in this decision was the following: a method is considered superior to another when the correlation between samples from the same speaker is closer to 1, and the cross-correlation (different speakers) is further from 1.)

In the Bravais-Pearson method, the LTS analysis spectrum is considered as a vector of dimension k , with a total of k frequency channels. The spectrum can then be defined as:

$$S \equiv (S_1, S_2, \dots, S_i, \dots, S_k), \quad (1)$$

where S_i is the level of the i -th frequency component (Harmegnies [12]). In this context, the coefficient R measures the relationship between

the two LTS samples. R is defined as:

$$R_{SS'} \equiv \frac{1}{k} \frac{\sum_{i=1}^k (S_i - M_S)(S'_i - M_{S'})}{\sigma_S \sigma_{S'}}, \quad (2)$$

where M_S and $M_{S'}$ refer to the averages of each spectrum, while σ_S and $\sigma_{S'}$ are the respective standard deviations. The Bravais-Pearson coefficient has several advantages, as it not only has a high discriminatory capacity, but is also independent of differences in relative intensity between the two spectra (Harmegnies [12]). This allows for comparing recordings that were made under different environmental conditions or microphone positioning.

Most studies on the relationship between speech and emotions have focused on the ability to distinguish between different emotional states of the speaker through the analysis of the acoustic signal (Williams and Stevens [4]; Fuller [14]; Scherer [7]; Johnstone [6]; Harnsberger *et al.* [15]). In particular, LTS analysis has been used to try to identify emotional states and depression (Pittam [16]), although human filtering has sometimes been used in the process (Banse and Scherer [5]).

Less common is research on how emotions or the deliberate intention to fake the voice affect FSI. Rodman and Powell [17] recommend and plan research to study the effects of masking in FSI, although they do not carry it out. More related to the present article, Hollien and Majewski [9] study the effect that stress (induced by electric shocks in subjects) has on the LTS analysis of the signal, but they have not found significant impacts for the purposes of conducting an adequate FSI. Except for the aforementioned, as far as the authors of this article are aware, there are no studies on the effect of emotions on FSILTS.

3. MATERIALS AND METHODS

This section will describe both the way in which the study population was defined and the method of executing the respective recordings.

3.1 Subject recordings

The problem of forensic speaker identification (FSI) has many variables. Consider, for example, the gender, geographical origin, and age of the subjects (Hertrich and Ziegelmayr [18]; Linville [19]; Hollien and Majewski [9]; Pittam [16]; Yüksel and Gündüz [20]). For this reason, subjects with similar characteristics were chosen. This approach was intended to foreground the effect of emotion above other variables, thereby reducing the overall complexity of the problem. The subjects met the following criteria: all were men aged 18–25, and their geographical origin was the Greater Metropolitan Area of the Central Valley of Costa Rica (this region is comprised of the four largest cities in the country, located in the central region of the country, which has the highest population density). All subjects attended primary school in this region.

3.2 Recording conditions

The recording conditions were standardized as much as possible in order to try to obtain accurate results. Each of the speech samples was recorded in a private place where the subjects felt comfortable. Additionally, they were asked to speak fluently and spontaneously, not recited or learned. High-quality smartphone microphones were used, such as those from iPhone and Samsung Galaxy.

Depending on the emotional state, there were two types of recording: Normal recordings, in which subjects were asked to talk about their daily life in order to have the least amount of emotional response possible; and angry recordings, in which subjects were asked to provoke a state of anger by describing a personal situation that made them angry. The interviewer left the subject alone in this part to avoid the subject from feeling inhibited.

The length of the samples was 45 and 60 seconds for normal and angry cases, respectively (the additional time for angry cases allows for a transition period).

4. METHODOLOGY

A total of 32 subjects (who met the previously mentioned selection criteria) were interviewed. This number was chosen to have a statistically significant sample (National Institute of Standards and Technology [21]) from the Greater Metropolitan Area of the Central Valley of Costa Rica. Each interview consisted of three recordings: For 16 of the 32 interviewees, the following recording order was implemented: normal-normal-angry. For the other 16 interviewees, the order normal-angry-normal was used instead. The motivation behind having two types of recording sequences is to reduce the effects of possible systematic errors due to the order of sampling.

The subjects were taken to a private place where they were read a set of instructions; each subject made their recording separately from the others. The interviewer explained to each subject that the experience was part of a research project of the University of Costa Rica and that the contents of the recordings would not be used or listened to. It was emphasized that only the acoustic properties of the samples were of interest and that the level of anger would be determined solely by their own self-assessment at the end of the recording process (this feature makes the process a blind experiment). In addition, the nature of the project was described only in very general terms so that the voice production would be as natural as possible.

For normal recordings, speakers were instructed to talk for 45 seconds about their life (e.g., their day, their pet, a recent event, etc.). For angry recordings, they were asked to talk for one minute about something that made them angry, for example, someone they did not get along with or an event that had irritated them greatly. Subjects were instructed not to fake the emotion (e.g., by forcing it by raising their voice). They were left alone with the recorder and told to speak when they were ready. Once the angry recording was completed, each participant was asked to rate their level of anger on a numerical scale from 1 to 5, where “1” meant “I was not able to get angry at all”, “3” meant “moderately

angry”, and “5” meant “furious”.

Note that when the first recording order (normal-normal-angry) was used, the subjects did not know about the anger part until the last recording, whereas in the other order (normal-angry-normal), the speaker already knew about this aspect of the research during their last recording. As mentioned earlier, both orders were used and averaged to reduce any possible systematic bias.

5. DATA PROCESSING

The data processing on which this research is based was developed by our group (Ortega-Rodríguez *et al.* [11]) studying FSI (without the anger component) and has been extensively tested and optimized to ensure the best identification results. By optimization, we refer to testing different types of time-windowing (e.g., rectangular, Gaussian, Welch, Hanning, Hamming, etc.) and different types of sampling rates to determine which works best according to the criterion mentioned in Section 2.

The data processing corresponding to this article can be summarized as follows. First, thirty-second segments are obtained from the original recordings. In the case of recordings with a state of anger, this procedure is particularly important to eliminate the transition part from a normal state to a state of anger. The audio processing software Audacity 2.1.0 (Audacity Team [22]) was then used to obtain the Fast Fourier Transform (FFT) for each recording. A Hanning window with a sampling of 4096 values was used. The use of this parameter is due to the fact that it showed the best results in our previous exploratory studies. Each FFT was saved as a text file for further processing. Figures 1 and 2 show samples of such spectra for normal and angry cases, respectively.

Subsequently, a C++ program calculated the Bravais-Pearson correlation coefficient for the samples. Finally, the average, *standard deviation* (SD) of the sample, and *standard error* of the mean (SE) for the correlation coefficient measurements were obtained. Both cases were

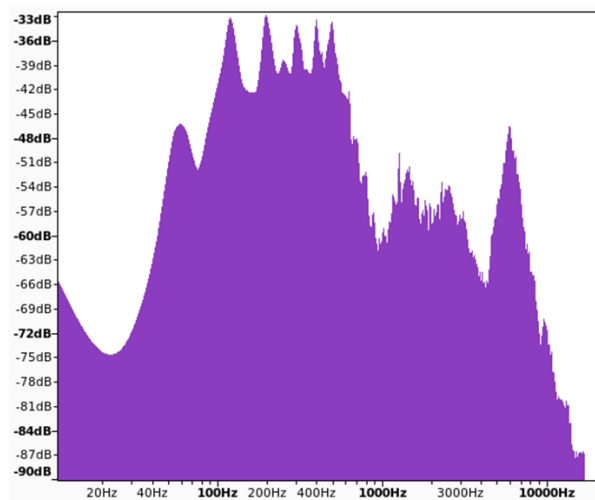


Figure 1: Spectrum of one of the recordings corresponding to normal speech (i.e., lacking anger). The sample duration is 30 seconds, and the spectrum was produced using a Hanning window with 4096 values.

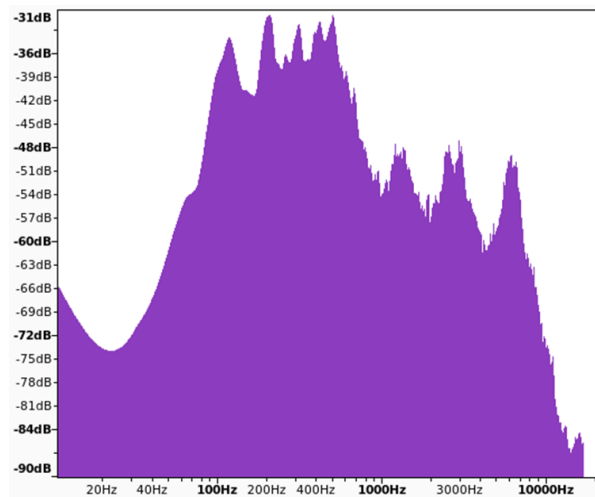


Figure 2: Spectrum of one of the recordings corresponding to level 4 anger speech (see Table 2). The sample duration is 30 seconds, and the spectrum was produced using a Hanning window with 4096 values. The speaker is the same as for the case of Figure 1.

performed: normal-normal-angry and normal-angry-normal.

6. RESULTS AND DISCUSSION

Table 1 shows the results of the calculations described in the previous section.

After the 32 subjects were interviewed, the Bravais-Pearson correlation coefficient R was calculated to obtain the correlation between normal and normal speech, and the correlation between normal and angry speech (intra-speaker).

Table 1: Statistical results of the Bravais-Pearson correlation coefficient R for the intra-speaker case. A total of 32 subjects participated in the contrast between the normal-normal and normal-angry cases.

	Correlation coefficient, normal-normal case	Correlation coefficient, normal-angry case
Average	0.950	0.934
Standard deviation	0.028	0.037
Standard error of the mean	0.005	0.005

As noted from the table, there is a noticeable difference between the two averaged values. To get an idea of how significant this difference is, it is useful to compare it with the results of our group’s aforementioned work (Ortega-Rodríguez *et al.* [11]), which is a simplified version of the process described in this article (since the element of anger was not previously present), although the experimental conditions were the same. In that work, the averaged correlation coefficient R between two different speakers was measured at 0.890 with an SE of 0.010, while the correlation between measurements of the same speaker had an average R of 0.955 with an SE of 0.005.

Comparing these three correlations: 0.950 (normal-normal, same speaker), 0.934 (normal-angry, same speaker), and 0.890 (different speakers, both in normal mode), it can be concluded that the effects of anger deviate the signal a significant 33% in the direction of a different speaker. This is very notable since the average self-reported anger was only 2.9 on the scale of 1 to 5 described earlier, which means that on average, the subjects were only moderately angry. The distribution of the anger level for the subjects can be seen in Table 2, and the respective statistical results are found in Table 3.

Table 2: Self-reported anger by the subjects; 1 means “I was not able to get angry at all”, 3 means “moderately angry”, and 5 means “furious”.

Number of subjects	Level of self-reported anger
5	4
18	3
9	2

It is expected that higher levels of anger would

Table 3: Statistical results for the data in Table 2.

	Self-reported anger level
Average	2.90
Standard deviation	0.65
Standard error of the mean	0.12

Table 4: Intra-speaker statistical results for the Bravais-Pearson correlation coefficient for the case of the five recordings with the highest anger. As expected, strong anger has a greater effect on FSILTS than moderate anger.

	Correlation coefficient, normal-normal case	Correlation coefficient, normal-angry case
Average	0.950	0.922
Standard deviation	0.030	0.047
Standard error of the mean	0.013	0.015

generate greater effects in FSILTS. To test this hypothesis, the results were filtered to only include the 5 most angry individuals (those with a level 4 on the scale of 1 to 5). Table 4 shows the results obtained under this condition. For strong anger, the deviation from normality is close to 50% of the way to another speaker.

This demonstrates that the effects of anger on FSILTS indeed grow as anger increases.

7. CONCLUSIONS

Although some authors praise the method of forensic long-term spectral analysis (FSILTS) for its robustness against speaker stress (Hollien and Majewski [9]), it has been found in this research that there is a significant distortion in the human voice due to anger for FSILTS purposes. Even when the emotional response of the participants remained moderate, a noticeable difference is found in the correlation coefficients between the cases of normal-normal and normal-angry recordings. Moderate anger deviates the results of FSI by 33% in the direction of another speaker. It should be emphasized that these results were obtained with a method that is fully automatable, providing an objective approach independent of human errors in perception. The method also avoids assessing the sincerity of the participants and is therefore in accordance with the code of practice of The In-

ternational Association for Forensic Phonetics and Acoustics (IAFPA [23]).

The results of this article are relevant to forensic research since LTS analysis has traditionally been considered a robust vector in FSI, especially because it is not sensitive to changes in speech sound intensity, works well for short recordings, and continues to function in the presence of noise and limited bandwidth. However, the results of this article indicate that care should be taken when using FSILTS to calculate likelihood ratios in the context of anger, even if this anger is not intense. In forensic applications, it is therefore recommended to always record the degree of anger of the speaking person, always keeping in mind the distortion values obtained in this article as a reference.

Since other emotions could also significantly affect the effectiveness of FSILTS, their study in automation processes is fully justified and recommended. The study could also be extended to women or speakers of other languages.

ACKNOWLEDGEMENTS

This work was supported by project 805-B2-175 of the Vice-Rectorate for Research of the University of Costa Rica, as well as the Geophysical Research Center of the same university.

AUTHOR CONTRIBUTIONS

All authors jointly worked on the generation of the general concept of this article and to define its methodology. They discussed and approved the results.

In particular, D. Valverde-Méndez and A. Venegas-Li were responsible for making the recordings, D. Valverde-Méndez conducted the data processing and wrote the first draft of the article (including background research), while M. Ortega-Rodríguez was responsible for the final text.

CONFLICT OF INTEREST

The authors declare no conflict of interest regarding the content of this article.

REFERENCES

1. HOLLIEN, Harry. Barriers to Progress in Speaker Identification with Comments on the Trayvon Martin Case. *Linguistic Evidence in Security, Law and Intelligence*, University Library System, University of Pittsburgh, v. 1, n. 1, p. 76–98, Dec. 2013. ISSN 2327-5596. doi: [10.5195/lesli.2013.3](https://doi.org/10.5195/lesli.2013.3).
2. HOLLIEN, Harry. An Approach to Speaker Identification. *Journal of Forensic Sciences*, Wiley, v. 61, n. 2, p. 334–344, Feb. 2016. doi: [10.1111/1556-4029.13034](https://doi.org/10.1111/1556-4029.13034), PMID: 27404606.
3. HOLLIEN, Harry Francis. *Forensic Voice Identification*. Londres, Inglaterra: Academic Press, 2002. ISBN 0123526213.
4. WILLIAMS, Carl E.; STEVENS, Kenneth N. Emotions and speech: Some acoustical correlates. *The Journal of the Acoustical Society of America*, Acoustical Society of America (ASA), v. 52, n. 4B, p. 1238–1250, Oct. 1972. doi: [10.1121/1.1913238](https://doi.org/10.1121/1.1913238).
5. BANSE, Rainer; SCHERER, Klaus R. Acoustic profiles in vocal emotion expression. *Journal of Personality and Social Psychology*, American Psychological Association (APA), v. 70, n. 3, p. 614–636, 1996. doi: [10.1037/0022-3514.70.3.614](https://doi.org/10.1037/0022-3514.70.3.614).
6. JOHNSTONE, Tom. *The effect of emotion on voice production and speech acoustics*. Dissertation (PhD) — University of Western Australia & University of Geneva, Perth, Australia, 2001.
7. SCHERER, Klaus R. Voice, Stress, and Emotion. In: _____. *Dynamics of Stress: Physiological, Psychological and Social Perspectives*. 1. ed. [N.p.]: Springer US, 1986. p. 157–179. ISBN 978-1-4684-5122-1. doi: [10.1007/978-1-4684-5122-1_9](https://doi.org/10.1007/978-1-4684-5122-1_9).
8. MARTIN, Maryanne. On the induction of mood. *Clinical Psychology Review*, Elsevier BV, v. 10, n. 6, p. 669–697, Jan. 1990. ISSN 1873-7811. doi: [10.1016/0272-7358\(90\)90075-1](https://doi.org/10.1016/0272-7358(90)90075-1).
9. HOLLIEN, Harry; MAJEWSKI, Wojciech. Speaker identification by long-term spectra un-

- der normal and distorted speech conditions. *The Journal of the Acoustical Society of America*, Acoustical Society of America (ASA), v. 62, n. 4, p. 975–980, Oct. 1977. ISSN 1520-8524. doi: [10.1121/1.381592](https://doi.org/10.1121/1.381592).
10. KINNUNEN, Tomi; HAUTAMAKI, Ville; FRANTI, Pasi. On the Use of Long-Term Average Spectrum in Automatic Speaker Recognition. In: *Proc. International Symposium on Chinese Spoken Language Processing*. [n.p.], 2006. p. 559–567. Available on <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=d3b4740466aeb1d25831b6329599b615a5bab9b1>.
11. ORTEGA-RODRIGUEZ, Manuel. *Final Report: Articulation of a speaker identification system for forensic purposes (Original: Informe Final: Articulación de un sistema de identificación de locutor con fines forenses)*. [N.p.], 2016. Accessed on November 2021. Available on <https://hdl.handle.net/10669/85190>.
12. HARMEGNIES, Bernard. SDDD: A new dissimilarity index for the comparison of speech spectra. *Pattern Recognition Letters*, Elsevier BV, v. 8, n. 3, p. 153–158, Oct. 1988. ISSN 1872-7344. doi: [10.1016/0167-8655\(88\)90093-1](https://doi.org/10.1016/0167-8655(88)90093-1).
13. STANTON, Jeffrey M. Galton, Pearson, and the Peas: A Brief History of Linear Regression for Statistics Instructors. *Journal of Statistics Education*, Informa UK Limited, v. 9, n. 3, Jan. 2001. ISSN 1069-1898. doi: [10.1080/10691898.2001.11910537](https://doi.org/10.1080/10691898.2001.11910537).
14. FULLER, Fred H. *Detection of emotional stress by voice analysis final report*. Bethesda, Maryland, USA, 1972. Available on <https://www.ojp.gov/ncjrs/virtual-library/abstracts/detection-emotional-stress-voice-analysis-final-report>.
15. HARNSBERGER, James D.; HOLLIEN, Harry; MARTIN, Camilo A.; HOLLIEN, Kevin A. Stress and Deception in Speech: Evaluating Layered Voice Analysis. *Journal of Forensic Sciences*, Wiley, v. 54, n. 3, p. 642–650, May 2009. ISSN 1556-4029. doi: [10.1111/j.1556-4029.2009.01026.x](https://doi.org/10.1111/j.1556-4029.2009.01026.x).
16. PITTAM, Jeffery. The Long-Term Spectral Measurement of Voice Quality as a Social and Personality Marker: A Review. *Language and Speech*, SAGE Publications, v. 30, n. 1, p. 1–12, Jan. 1987. ISSN 1756-6053. doi: [10.1177/002383098703000101](https://doi.org/10.1177/002383098703000101).
17. RODMAN, Robert D.; POWELL, Michael S. Computer Recognition of Speakers Who Disguise Their Voice. In: *The International Conference on Signal Processing Applications and Technology (ICSPAT 2000)*. [n.p.], 2000. Available on <https://api.semanticscholar.org/CorpusID:16980245>.
18. HERTRICH, I.; ZIEGELMAYER, G. Sexual dimorphism in the long term speech spectrum. *Human Evolution*, Springer Science and Business Media LLC, v. 2, n. 3, p. 255–262, May 1987. doi: [10.1007/bf03016110](https://doi.org/10.1007/bf03016110).
19. LINVILLE, Sue Ellen. Source Characteristics of Aged Voice Assessed from Long-Term Average Spectra. *Journal of Voice*, Elsevier BV, v. 16, n. 4, p. 472–479, Dec. 2002. doi: [10.1016/s0892-1997\(02\)00122-4](https://doi.org/10.1016/s0892-1997(02)00122-4).
20. YÜKSEL, Mustafa; GÜNDÜZ, Bülent. Long term average speech spectra of Turkish. *Logopedics Phoniatrics Vocology*, Informa UK Limited, v. 43, n. 3, p. 101–105, Sep. 2017. doi: [10.1080/14015439.2017.1377286](https://doi.org/10.1080/14015439.2017.1377286).
21. National Institute of Standards and Technology. *NIST/SEMATECH e-Handbook of Statistical Methods*. [n.p.], 2012. Accessed on October 2021. Available on <https://www.itl.nist.gov/div898/handbook/prc/section2/prc222.htm>.
22. Audacity Team. *Audacity (v. 2.1.0), audio editor and recorder*. 2015. Available on <https://www.audacityteam.org/>.
23. The International Association for Forensic Phonetics and Acoustics. *Code of Practice*. [N.p.], 2004. Accessed on January 2018. Available on <https://www.iafpa.net/the-association/code-of-practice/>.