

On the design of automatic voice condition analysis systems. Part I: review of concepts and an insight to the state of the art.

J.A. Gómez-García^{a,*}, L. Moro-Velázquez^{b,a}, J.I. Godino-Llorente^a

^aUniversidad Politécnica de Madrid. Ctra. Valencia, km. 7, 28031. Madrid, Spain.

^bJohns Hopkins University. Baltimore, Maryland 21218, USA.

Abstract

This is the first of a two-part series devoted to review the current state of the art of automatic voice condition analysis systems. The goal of this paper is to provide to the scientific community and to newly comers to the field of automatic voice condition analysis a resource that presents introductory concepts, a categorisation of different aspects of voice pathology and a systematic literature review that describes the methodologies and methods that are mostly used in these systems. To this end, the phenomena of pathological voice is firstly described in terms of perceptual characteristics and its relationship with physiological phenomena. Then, a prototypical automatic voice condition analysis system is described, discussing each one of its constituting parts and presenting an in-depth literature review about the methodologies that are typically employed. Finally, a discussion about some variability factors that affect the performance of these systems is presented.

Keywords: Automatic voice condition analysis, voice pathology detection, extralinguistic aspects of the speech, voice quality.

1. Introduction

Speech is accomplished through complex articulatory movements that mould the vocal excitation source in order to convey spoken sounds. In this process, three components can be identified: The *excitation source* (be it voiced, unvoiced, a mixture of both or its absence -such as in a pause-) providing the driving force for the speech production process, the *articulation* defined by the movements of the speech articulators moulding the production of a certain sound, and the *fluency* defining the rate at which the speech is generated. Despite the main objective of speech is transmitting information by means of sounds that encode linguistic content, the inherent intricacy of the production process embeds a substantial amount of non-linguistic information that is often described in terms of dimensions [1–3]. In this regard, the *paralinguistic* dimension of speech conveys information about the affective, attitudinal or emotional state of the speaker; the *extralinguistic* dimension informs about the speaker’s identity and state (with *traits* such as age, sex, condition, etc.); whereas the *linguistic* dimension is related to the message, variations in language, dialect, sociolect, idiolect and speech style of the speaker. There is often described a fourth *transmittal* dimension that tells nothing about the speaker but about its physical location.

*Corresponding author

32 By virtue of the valuable information contained in speech and the inexpensiveness and easiness of capturing
33 ing speech signals in a non-invasive manner, a great deal of interest has arisen in designing systems capable
34 of isolating a certain dimension (or trait within a dimension) for automatic analysis purposes. As a matter
35 of example, the literature presents systems that have been focused on extracting extralinguistic information
36 to automatically determine the identity of speakers [4], age [5] or sex [6]. In the same manner, paralinguis-
37 tic information has been extracted to identify the speaker's emotions [7] or level of interest [8]. Likewise,
38 the linguistic dimension has been studied to recognize the accent, dialect [9] or the message itself (speech
39 recognition) [10].

40 One application that has been gaining popularity during the last years is in the *analysis of speakers'*
41 *condition* using voice recordings. In this respect, the clinical evaluation of voice disorders often relies on an
42 instrumental examination and a perceptual analysis of the speech. The *instrumental* medical examination
43 focuses on a primary aetiological diagnosis through the investigation of acoustic, aerodynamic, electroglotto-
44 graphic, videolaryngostroboscopic and/or the exploration of other types of biosignals; whereas the *perceptual*
45 examination extracts multidimensional information that is not quantifiable instrumentally, by means of a
46 qualitative description of the perceived degree of dysphonia that is present in the voice [11–13]. This infor-
47 mation might be complemented by an interview where the patient states his/her symptoms, the examination
48 of his/her medical records, evaluation of other body functions and systems, and exploration of the laryngeal
49 structures and their function. These procedures should lead the medical expert to a diagnosis about the
50 condition of the patient. The diagnostic process is differential, i. e., all possible causes of a problem are
51 considered, and then the available information is matched against each one of the hypothesis explaining the
52 disorder in the search for a match [14]. In some severe pathological cases a decision about the condition of
53 the patient is straightforward, but in others, it would probably be conditioned by subjective factors or the
54 observations and hypothesis made by the clinician. The increasing need of improving the diagnosis of voice
55 pathologies has given rise to an emerging field called *Automatic Voice Condition Analysis (AVCA)*, that
56 aims at analysing, classifying and quantifying the degree to which a patient is affected by a voice disorder.
57 This analysis is performed using automatic systems that provide objective measurements of the patient's
58 vocal condition, exploiting the close relationship that exists between acoustic features extracted from the
59 speech and voice pathology [15]. This reduces the evaluation time and the cost of diagnosis and treatment,
60 providing extra advantages such as the avoidance of invasive procedures thanks to the employment of speech
61 signals which are easily recorded by inexpensive means [16].

62 With these precedents in mind, the aim of this paper is to provide a review of AVCA systems, introducing
63 key concepts related to vocal pathologies and their acoustic consequences in voice signals. A typical AVCA
64 system is also described, detailing each one of its constituting blocks while providing a thorough literature
65 review to portray the most used methodologies. Finally, some confounding factors that affect AVCA systems
66 are discussed as well. It is worth noting that the main interest of this paper is related to the automatic *dete-*
67 *ction* (classification of control vs. pathological) and *identification* of voice disorders (multiclass classification

68 of the actual disorder affecting speech), rather than the *assessment or grading* of voice signals. Indeed, we
69 consider that the assessment of voice pathologies deserves of a separate paper to handle the particularities
70 of this classification task.

71 This paper is organised as follow: section 2 describes introductory concepts related to voice pathologies;
72 section 3 introduces AVCA systems, whereas section 4 describes its constituting blocks, presenting a review
73 of the most typical methodologies employed in literature. Then, section 5 describes some variability factors
74 affecting this type of systems. Finally, section 6 presents some discussions and concluding remarks.

75 2. Voice pathologies

76 Following the description presented in the introduction, a *speech disorder* can be defined as an impairment
77 of the articulation of speech sounds, fluency and/or voice [17, 18]. It is worth noting that from all these
78 elements, *this paper is only focused in those pathologies affecting voice*, and therefore the main interest of
79 this paper is in the study of phonatory aspects of the speech. Articulatory, prosodic or language disorders
80 are by themselves topics that should be handled separately.

81 To address the concepts of *voice condition* it is firstly necessary to describe the properties of a "normal"
82 voice. This, however, poses numerous difficulties since there exist several definitions of "normality", and
83 the distinction from what can be considered healthy or abnormal relies on subjective perceptual judgements
84 made by listeners or by the speaker itself [19]. Indeed, a singer who uses a deviant voice as a trademark,
85 might acknowledge his/her voice as normal, but this can be perceived otherwise by some listeners. By
86 contrast, a high-pitched voice which in different circumstances would be considered normal if uttered by a
87 child, might be deemed as pathological if uttered by an adult. In spite of that, there are certain common
88 characteristics that can be regarded as normative, and thus, can be utilised as synonyms of *non-pathological*
89 *voice condition*. Literature presents definitions which differ in terms of what (and how) can be categorised
90 as normal or **normophonic**. In this paper, we adopt the *perceptual* definition presented in [20], on which a
91 normophonic voice is described in terms of the following properties: (i) a pleasant quality, with an absence
92 of noise, inappropriate breaks, perturbations or atonality; (ii) pitch in accordance to the age and sex of the
93 speaker; (iii) loudness that is appropriate to the communication event; (iv) pitch and loudness variations
94 that are available to express emphasis, meaning or subtleties indicating individual feelings and semantic
95 differences; (v) sustainability to meet social and occupational needs.

96 Abnormal voices do not possess any, a combination, or all of the above properties. Typically, three types
97 of aberrant voices are usually identified: aphonia, dysphonia and muteness [20]. **Aphonia** is characterised
98 by the absence of vibration of the vocal folds -but not of sound- resulting in a voice that is perceptually
99 described as extremely breathy. Similarly, **muteness** is referred to the absence of vocal folds vibration,
100 accompanied by the inability to produce audible sounds. Finally, **dysphonia** is described by the absence of
101 vocal quality, pitch, loudness, and/or variability which is inappropriate for an individual's age and/or sex
102 [17, 20].

103 From the perspective of AVCA systems, muteness has to be discarded from the study due to the unavail-
104 ability of audible outputs for automatic analysis purposes using voice registers. Similarly, the perceptual
105 consequences of aphonia are so notorious that an automatic analysis to detect or assess the impairment is
106 seldom considered. Consequently, and to the authors' knowledge, there is not a single work in literature
107 dealing with automatic analysis of aphonic voices. At the end, only dysphonic and normophonic voices are
108 examined by AVCA systems for labours of identification, detection or grading of pathological states.

109 Revisiting the definitions of dysphonia and normophonia it can be observed that 4 elements are identified,
110 i. e., loudness, pitch, quality and variability. In this regard, *loudness* is defined as a perceptual correlate of the
111 intensity of the sound pressure created by the release of air through the glottis. Disorders affecting loudness
112 occur when the voice is louder or softer in concordance to the speaker's context. Loudness impairments are
113 often indicators of personality disorders (overly aggressive, shy, or socially insecure behaviour), or are the
114 consequence of certain pathologies such as Parkinson's disease or paresis. The second perceptual trait is
115 *pitch*, which is the correlate of the frequency of vibration of the vocal folds, i. e., the fundamental frequency
116 (f_0). The rate of vibration is determined by the physical characteristics of the vocal folds such as the mass,
117 elasticity or length. Impairments affecting pitch include those where voice is tremulous or abnormal in
118 concordance to the speakers' context. Examples include mutational falsetto (abnormally high-pitched voice
119 uttered by an adult not correlating with his age or sex) or ventricular phonation (abnormally low-pitched
120 voice product of the vibration of the false vocal folds). The third perceptual trait is *quality*, which is a
121 correlate of the vibrational patterns of the vocal folds and resonant characteristics of the vocal tract. This
122 is in turn composed of other traits describing very differentiated physiological phenomena as described next
123 [18, 20]: (i) *strain*, related to disturbances in the vibratory patterns of the vocal folds due to an excessive
124 tension in the larynx which may result in over adduction of vocal folds; (ii) *breathiness*, related to turbulent
125 air streams released during incomplete vocal closure; (iii) *roughness (or harshness)*, related to irregularities
126 or vibration defects of the vocal folds; and (iv) *resonance*, caused by abnormalities present in the vocal
127 tract, such as defects in the closure of the velopharyngeal port. Despite there is a fourth trait related to
128 *variability*, examining the flexibility of voice in relation to variations of pitch, loudness and quality in spoken
129 contexts, we consider that -at least in AVCA methodologies- these variations can be included directly into
130 their respective descriptors. We also believe that is possible to consider a superclass that embeds pitch and
131 loudness into a sphere of *vocal aspects* as both examine the concordance of the voice respecting the speaker
132 or its context. Therefore, we propose to analyse voices using automatic systems on the basis of two spheres,
133 one related to the *vocal aspects* on which the adequacy of the pitch, loudness and its variations are examined
134 in relation to the context and to normative values of an average speaker in the same population group. And
135 a sphere related to *vocal quality* on which the resonant characteristics and the vibrational patterns of the
136 vocal folds are examined.

137 The sphere of vocal aspects has been seldom studied in literature, with a vast majority of papers published
138 in topics related to the analysis of vocal quality. In addition, and despite resonant patterns make up for

139 voice quality, the most important descriptor often found in the state of the art is referred to a quantification
 140 of the *vocal aperiodicity* describing the manner of vibration of the vocal folds. Due to their importance in
 141 the design of AVCA systems, the mechanisms of vocal aperiodicity are described next.

142 It has been stated that three¹ types of vibrational patterns are often encountered in voice signals [22]:
 143 (i) *type I* voices, characterised by a nearly periodic behaviour; (ii) *type II* signals, which contain bifurcations,
 144 sub-harmonics or modulating frequencies; and (iii) *type III* voices, which are characterised by an aperiodic
 145 behaviour. In accordance with such distinction, normophonic voices are usually enclosed into the Type I
 146 typology, whereas pathological voices are embodied into the Type II and III categories [22, 23]. As a matter
 of example, Figure 1 illustrates some cases of voice signals following the above-mentioned typology.

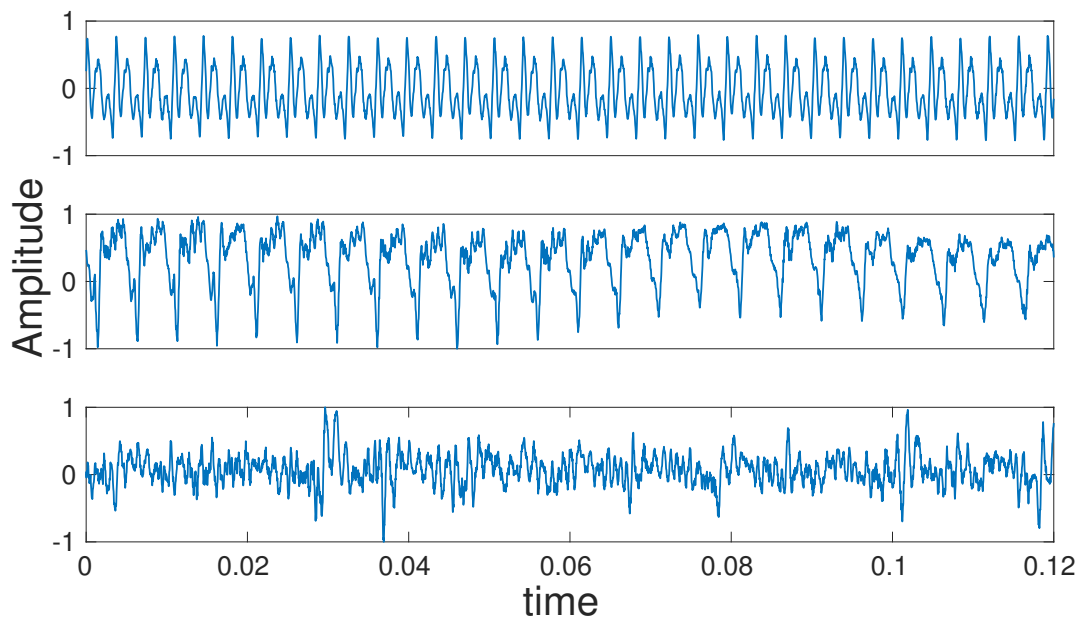


Figure 1. Typology of voice signals according to [22]. *Top panel:* normophonic type I signal, characterised by a periodic behaviour; *middle panel:* pathological type II signal having modulating frequencies; *bottom panel:* pathological type III signal characterised by an aperiodic behaviour.

147
 148 Vocal aperiodicity is explained as the result of some very distinctive processes occurring during the voice
 149 production process such as [24–27]: (i) *irregular dynamics* of the vocal folds and involuntary transients
 150 between dynamic regimes (distinguishing features on very specific voice impairments such as diplophonia
 151 or biphonation); (ii) *modulation noise* owing to extrinsic perturbations in amplitude and frequency of the
 152 glottal cycle and which is often associated to roughness; and (iii) *additive noise* owing to turbulent airflow
 153 and which is correlated to breathy vocal quality.

154 To summarise some of the concepts introduced in this section, Figure 2 is presented.

¹Despite there are authors defining a fourth category [21], a vast majority of works still employ the most classical definition.

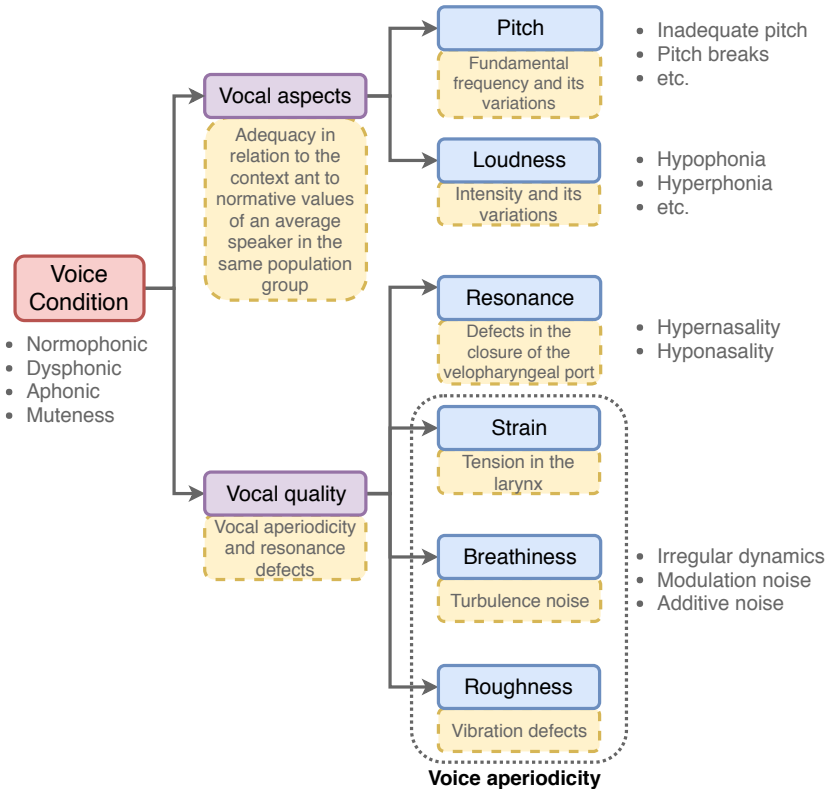


Figure 2. Graphical representation of the concepts introduced in the section. The dashed lines are referred to the physical phenomena whereas the box in blue is referred to the perceptual trait.

155 **3. Automatic voice condition analysis systems**

156 Generally, AVCA systems follow a pattern recognition-like structure on which characteristics are extracted
 157 from the acoustic signal in the form of a set of features to accomplish a further decision making task. An
 158 example of a typical AVCA scheme is presented in Figure 3.

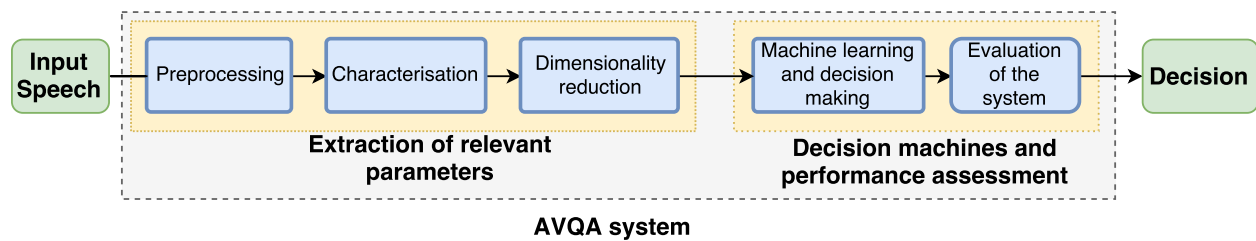


Figure 3. Depiction of a typical AVCA system.

159 Before going deeper into each one of the building stages of AVCA systems, two initial considerations
 160 -referred to the *input speech* and *decision* blocks in the depiction- are to be addressed first in the following
 161 subsections to respond to the questions: (i) what type of speech task is to be used for the design of the
 162 system; (ii) what type of decision should the system provide.

163 3.1. *The input speech*

164 The interest of considering the acoustic material in the design of AVCA systems arises from the fact
165 that depending on the type of utterance, different configurations of the speech production subsystems arise,
166 permitting the analysis of certain aspects of speech or others. Indeed, some pathologies are more likely to be
167 identified when examining determined units of speech. For instance, resonance impairments are more easily
168 perceived in utterances containing /m/ or /n/ prompts.

169 In general, two types of *speech production tasks* are employed for the evaluation of voice condition:
170 sustained phonation of vowels and running speech. On one hand, *sustained phonation* is the result of the
171 production of voiced sounds due to the vibration of the vocal folds, as when a vowel is uttered and maintained
172 during a certain amount of time. Some advantages of using sustained phonations in AVCA systems include
173 [23, 28]: (i) the facility to be analysed by automatic tools; (ii) the production of vowels is straightforward;
174 (iii) vowels are not affected by paralinguistic or extralinguistic characteristics such as speaking rate, dialect,
175 intonation, and idiosyncratic articulatory behaviour; (iv) vowels often generate simpler acoustic structures
176 that might lead to consistent and reliable perceptual judgements of voice quality; (v) vowels do not depend
177 on extra processing stages (such as voiced/unvoiced detectors) for the design of AVCA systems.

178 The selection of the vowel to be uttered is also a relevant matter. It has been stated that the type of
179 vowel -along with the vocal effort and the muscle tension in the larynx- influences the degree of vocal folds
180 approximation, affecting the perception of voice quality [29]. For this reason certain open vowels, such as
181 /a/, are often employed in AVCA systems since they are produced with a relatively open tract allowing the
182 examination of the entire vocal tract apparatus. By contrast vowels like /i/ and /u/, may not allow this
183 examination due to the separation between the front and back cavities of the mouth during its production
184 [30].

185 On the other hand, *running (or connected) speech* is the result of the source signal (either voiced, unvoiced,
186 a mixture of both or its absence) being modulated by the articulatory subsystems, as when uttering a certain
187 word or pronouncing sentences. Despite this speech production task is not as widely popular in AVCA
188 systems as the one based on sustained vowels, there exist strong arguments favouring its use. Indeed, one
189 interesting property of running speech comes from a phenomenon called *coarticulation*, which is related to the
190 influence of the preceding and succeeding acoustic unit on the current unit under analysis. The dynamical
191 effects introduced by coarticulation might be relevant for certain applications. Besides that, it has been
192 stated that the impressions of certain characteristics of vocal quality are more easily perceived on vowels
193 generated in a voiced context, vowels after a glottal closure, or during the production of strained vowels [31].
194 Some additional advantages reported in literature in favour of analysing connected speech include [15, 25]:
195 (i) it requires switching on and off the vibration of the vocal folds continuously, or maintaining the voicing
196 while the supraglottal apparatus changes, facilitating the exploration of certain dynamic aspects of the
197 speech production; (ii) speakers are less likely to compensate for voice problems while producing connected
198 speech than while phonating sounds; (iii) running speech provides a more realistic scenario since sustained

199 phonations are more characteristic of singing rather than speaking; (iv) running speech contains fluctuations
200 of vocal characteristics in relation to voice onsets, terminations and breaks; (v) in certain voice disorders
201 (e.g., spasmodic dysphonia), the production of sustained vowels is less symptomatic than in connected
202 speech, which may lead to an underestimation of the impairment; (vi) running speech contains variations in
203 pitch and loudness, parameters that are important in the analysis of abnormal voice quality.

204 3.2. Automatic decision tasks

205 Three fundamental tasks may be considered in an AVCA system: voice pathology detection, voice pathol-
206 ogy identification and voice pathology assessment.

207 On one hand, *voice pathology detection* is a two-classes decision making process aiming to decide whether
208 a given speech register is normophonic or pathological (dysphonic or aphonic). On the other hand, *voice*
209 *pathology identification* is a multi-class decision making process on which the goal is to assign a category
210 to the input speech. The identification task is typically made in terms of the actual pathology (nodules,
211 Reinke’s oedema, etc.), the aetiology (organic, functional, etc.) or any other categorisation that groups
212 general aspects of the analysed speech. From a practical point of view, identification is more challenging than
213 detection, because of the multiple-class scenario on which it is defined. Nonetheless, both tasks are intricately
214 complex due to several factors, such as the wide range of profiles that are found for normophonic voices, the
215 documented overlap between normophonic and pathological states [32], the close relationship between disease
216 and certain quality factors associated to normal processes such as ageing [33], the simultaneous presence of
217 pathologies of different aetiologies in the same patient, etc.

218 By contrast, *voice pathology assessment* is aimed at grading the level of pathology that is perceived in a
219 given speech signal. This is of great relevance since it is not possible to instrumentally delimit a phonation
220 behaviour categorically. An useful descriptor of dysphonia is the **hoarseness**, which portrays the noisy,
221 atonal and/or odd vocal resonance patterns encountered in voices [20]. The hoarseness is widely employed
222 in literature, as perceptually, it is described as a superclass that contains roughness and breathiness -
223 the two most reliable traits describing vocal quality- [31]. The assessment task is generally performed
224 in concordance to a *perceptual rating scale* that evaluates voice quality and provides information about
225 the level of impairment. Some popular perceptual evaluation scales include the GRBAS [12, 31], Voice
226 Handicap Index [34] and CAPE-V [35, 36] scales. Even though the perceptual scales have been designed
227 to evaluate every aspect that is relevant to voice quality, the reliability of the ratings is conditioned by the
228 multidimensional aspects of voice quality, the intrinsic variability of speech, the subjectivity of perception
229 [37], and the nonlinear relationship between pathology and measured or perceived voice quality [32].

230 4. Prototypical AVCA systems: an insight to the state of the art

231 The present section describes each one of the building blocks of Figure 3, providing a review of some
232 relevant techniques often employed. Along with the description, a literature review is provided to identify

233 the techniques which have been used before by other authors.

234 The literature review is performed using the web search of Scopus[®] and Google scholar[®]. The terms that
235 are employed include: "dysphonia", "pathology", "automatic", "voice", "quality", "classification", "detection",
236 "identification" and combinations and derivations of them. Pathological states including "Parkinson's
237 disease", "Alzheimer", "Obstructive Sleep Apnoea", "Nodules", etc. are also used in the web search. The
238 review is limited to those papers published after the year 2000, focusing on journal papers (although some
239 documents in reputable conferences are also included) listed in the Journal Citation Reports[®] or the Scimago
240 Journal Rank[®]. The list of predatory journals (<https://predatoryjournals.com/journals/>) is also con-
241 sulted to discard papers published in journals engaged in predatory practices.

242 4.1. Input speech

243 The collection of exemplar recordings describing the classes under study conforms a *dataset* or *corpus* of
244 speakers which are typically used to train and test an AVCA system. Although the term *database* is often
245 found in the literature, we discourage its use to avoid the technical connotations that it has in computer
246 science.

247 The data acquisition process should follow certain guidelines to prevent the introduction of unexpected
248 variability, including the avoidance of external sources of noise or the preservation of similar acoustic and
249 instrumental conditions during the recording process. Some recommendations for the acquisition of voice
250 signals for acoustic voice analysis has been presented in [38], advising -among others- the use of professional
251 condenser microphones with a minimum sensitivity of -60 dB, constant mouth-to-microphone distances less
252 than 10 cm, sampling frequencies between 20 to 100 kHz, and sound-treated rooms with ambient noise lower
253 to 50 dB. Other considerations in terms of technical characteristics of microphones have been described in
254 [39], where a flat frequency response microphone is recommended, within the frequency of lowest fundamental
255 frequency and highest spectral component, equivalent noise level at least 15 dB lower than the sound level
256 of the softest phonation, etc. In the same way, the recorded corpus should be large enough to contain all
257 possible variations within the class, while being balanced in terms of age, sex, etc. , properly representing
258 variations in speech due to accent, dialect, sociolect, etc. Likewise, conditions such as smoking or professional
259 voice should be accounted. The management of the corpus for the storage and accessibility of recordings from
260 a medical perspective should also be considered, as this permits the creation of synergies towards certain
261 tasks such as the diagnosis of the pathologies, or the assessment from a perceptual point of view as given by
262 different evaluators. Some considerations referred to the management in a clinical setting for a large corpus
263 of dysphonic and dysarthric speakers are discussed in [40].

264 Literature reports the existence of several public and private datasets that have been recorded for
265 the purposes of detection, identification or assessment of voice pathologies. Regarding *public* datasets, the
266 *Massachusetts Ear and Eye Infirmary* (MEEI) [41] is probably the most widely employed corpus, being for
267 years, the sole resource that was available for the study of pathological speech. MEEI contains approximately
268 700 registers of the vowel /a/ and the first sentence of "*the rainbow passage*" text, recorded at varying

269 sampling frequencies (25 kHz to 50 kHz). A subset of the corpus -chosen to ensure a balance in age, sex and
270 pathologies- has been published in [42], becoming a standard partition for comparisons among different works
271 in literature. Despite its popularity, MEEI suffers from well-known problems which might bias outcomes, like
272 the different recording conditions of normophonic and pathological speakers [43, 44]. Recently, another speech
273 pathology corpus has been made accessible publicly: the SVD dataset [45]. This partition was recorded by
274 the Institut für Phonetik at Saarland University and the Phoniatriy Section of the Caritas Clinic St. Theresia
275 in Saarbrücken, Germany. It contains more than 2000 acoustic and *electroglottographic* (EGG) registers of
276 the vowels /a/, /i/ and /u/ phonated at normal, high, low, and rising-falling pitch; as well as registers of
277 the sentence "*Guten Morgen, wie geht es Ihnen?*" (*Good morning, how are you?*), recorded at 50 kHz and
278 16 bits of resolution.

279 Some of the most well-known *privative* datasets include the *Hospital Príncipe de Asturias* (HUPA) corpus
280 [46] which contains registers of the sustained vowel /a/ of 366 adult Spanish speakers (169 pathological and
281 197 normophonic); or the *Arabic Voice pathology Dataset* (AVD) [47, 48] which is composed of registers of the
282 vowel/a/ and running speech of 188 normophonic and 178 pathological Arabic speakers. Another notable
283 corpora, which is perhaps one of the largest in terms of number of patients, is the one recorded in hospitals
284 in Marseilles and Aix-en-Provence in France [40]. It is composed by registers of *sound-pressure level* (SPL),
285 oral airflow, and subglottal air pressure of more than 2500 dysphonic, dysarthric and normophonic speakers.
286 There exist other *privative* corpora, exhibiting a large variety of characteristics respecting the acoustic
287 conditions followed during the recording process, the instrumentation, the type of speech material that is
288 elicited, the type of disorders, etc. Indeed, most of the *privative* datasets contain microphonic recordings
289 of the sustained phonation of vowel /a/ [49–67] or a combination of several vowels [68–72]. There are some
290 datasets with registers of running speech for different languages, which in text-dependent scenarios employ
291 isolated words [73, 74], reading of phrases such as "*the rainbow passage*" [75–77], "*the north wind and the*
292 *sun*" [78, 79], "*the story of Arthur the rat*" [80], or other texts [73, 81–86]. In the text-independent case
293 they employ conversational speech or other types of elicitation tasks [84]. There exist some other datasets
294 that include other type of complementary biosignals besides the acoustic recording. Namely, some contain
295 EGG recordings of the vowel /i/ [87, 88], /a/ [51, 66]; while others complement the acoustic registers with
296 questionnaire data [89, 90] or laryngoscopy information [90].

297 Regarding the languages that are reported, literature indicates datasets uttered by Russian [57, 74],
298 Korean [91], Spanish [46], Colombian [92], Arabic [48], German [32, 93], Czech [94], Dutch [85, 86], Chinese
299 [84], Brazilian [69], French [40] or Lithuanian [95, 96] speakers. Likewise, some include a broad range of
300 voice pathologies [41, 45, 46], whereas other are concentrated in certain disorders such as nodules [97],
301 polyps [59, 98], larynx cancer [87, 88, 99, 100], hypofunctional voices [11], diplophonia [101], spasmodic
302 disphonia and muscle tension dysphonia [102], unilateral laryngeal paralysis [50, 103], obstructive sleep
303 apnea [104, 105], hypernasality [106], Parkinson’s disease [92, 107, 108], dysphagia [109, 110], lupus [111],
304 etc.

305 4.2. Preprocessing

306 Since speech is intrinsically non-stationary, some preprocessing methods are often employed before the
307 utilisation of conventional signal processing techniques which rely on stationary requirements. One common
308 procedure is the *short-time analysis*, which decomposes the input signal into a series of equal-length chunks of
309 speech, called *frames*, permitting the treatment of each individual chunk as a stationary or quasi-stationary
310 fragment. This procedure is composed of two operations: *framing*, which divides the signal into frames
311 (typically overlapped); and *windowing*, which tapers the beginning and ending of the frames, through the
312 product with a window function to improve spectral properties. The window function should be selected to
313 provide a frequency response with a narrow bandwidth in the main lobe and large attenuation in the side-
314 lobes. Popular choices include the triangular, Hanning or Hamming windows, whereas the window length
315 varies depending on the application. Typically for applications using sustained phonation, the duration of
316 the window is set in between 20-40 ms. The upper limit ensures that frames are not that large to make the
317 quasi-stationarity assumption void, whereas the lower limit is set to make the analysis independent of the
318 location of pitch pulses within the segment, while ensuring at least two to three times pitch periods (since
319 the typical range of pitch frequency is between 80 – 500 Hz, a pitch pulse is expected every 12 – 2 ms [112]).
320 For applications using running speech, window lengths are typically set in the order of 20 – 30 ms to conserve
321 the quasi-stationary assumptions [113].

322 From those works reporting the type of window that is employed, literature indicates the popularity
323 of Hamming [11, 23, 28, 37, 43, 44, 57, 63, 74, 74, 81, 89, 93, 94, 114–120, 120–132, 132, 133, 133–143]
324 or Hanning [27, 88, 91, 96, 138, 144–150] windows. The length of the window varies depending on the
325 application, type of speech task and characteristics that are utilised. Popular values found in literature
326 include 10 ms [65, 91, 139, 144], 16 ms [151], 20 ms [48, 81, 121, 131–133, 152–154], 25 ms [74, 128, 155], 30
327 ms [68, 114, 120, 130, 143], 40 ms [43, 129, 134, 142, 146, 147, 156–163], 50 ms [61, 77, 88, 96, 118, 149] or
328 55 ms [141, 156]. Notwithstanding, for certain types of features they can be as large as 80 ms [164], 100 ms
329 [165], 200 ms, [23, 166], 262 ms [119], 400 ms [167], or 800 ms [168].

330 Other preprocessing techniques often found are voice/unvoiced and endpoint detectors [37, 43, 66, 146,
331 147, 162], which ensure that only segments generated during the vibration of the vocal folds are employed; or
332 silence detectors [169] that eliminate utterances not containing speech. Similarly, and to remove the influences
333 of the vocal tract in the speech signal, inverse filtering is often employed [52, 69, 120, 170–174]. Likewise,
334 the use of pre-emphasis filtering to accentuate the high frequency content of speech [74, 89, 121, 130] has
335 been employed, although it has been reported that it does not improve detection results in AVCA systems
336 [175].

337 4.3. Characterisation

338 The characterisation stage has the goal of extracting features capable of portraying the properties of the
339 classes under analysis. The idea is to extract a d -dimensional vector of characteristics, $\vec{x} = \{x[1], \dots, x[d]\}$,

340 describing d properties of the input speech. Usually this vector is associated to a label ℓ indicating the
341 membership of the utterance to a certain class, although this is not necessary. The features can be extracted
342 either in a *short-time basis* (as introduced in the previous section) having as many vectors of features as
343 frames of speech, or in a *long-term basis* calculating a single vector of characteristics per audio register.

344 Finding characteristics that effectively describe the presence of voice impairment is difficult, specially since
345 some phenomena associated to voice disorders (such as aperiodicity) are present even in non-pathological
346 states due to perturbations inherent to the phonation process [176]. As a result, there is no single fea-
347 ture, in the context of screening, that perfectly differentiates between normophonia and pathology, or that
348 biunivocally correlates acoustic measurements and voice quality [157]. A common approach to counteract
349 this, consists on studying different types of features, in the hope of finding combinations of characteristics
350 that complement with each other. In this sense, the best characteristics would be those with the lowest
351 correlation with the others but capable to provide the best discrimination capabilities [138]. Although mul-
352 tidimensional studies have reported good performance in screening tasks, this type of analysis is usually
353 carried out by complex pattern recognition techniques, which makes difficult the interpretation of results
354 from the perspective of a human evaluator [138].

355 The vast majority of descriptors of voice condition seek to compute metrics of vocal quality due to their
356 close relationship to features extracted from voice signals. Literature reports diverse characterisation schemes
357 which have been found to perform differently according to the pathologies under study or the dataset that
358 has been employed. In general, some popular features -to be described in the following sections- include
359 those based on *temporal and acoustic analysis*, *perturbation and fluctuation*, *spectral-cepstral*, *complexity*,
360 *3-dimensional representations* and *other types of features* not fitting in the above categories.

361 4.3.1. Temporal and acoustical analysis

362 To the best of the authors' knowledge, there are only a few papers accounting for the *vocal aspects*
363 of speakers (see Figure 2). In the first two, authors seek to monitor aberrant patterns of f_0 and SPL of
364 hyperfunctional speakers and employ them for the distinction between dysphonic and normophonic speakers
365 [77, 177]. Likewise, in [102], the degree of voice breaks are used to distinguish spasmodic dysphonia speakers
366 from others suffering from muscle tension dysphonia and a control population. There are not automatic
367 systems that seek to correlate SPL and f_0 to the contextual and personal traits of the speaker under analysis.

368 By contrast, most of literature is referred to the analysis of descriptors of vocal quality. With regards
369 to the analysis of *irregular dynamics* in pathologies such as diplophonia, literature reports the computation
370 of the degree of subharmonics and the diplophonia diagram [101, 178]. To capture *modulation* and *additive*
371 *noise*, some approaches are based on tracking f_0 and deriving low-order statistics to track disturbances in
372 the normal vibration patterns of the vocal folds [32, 59, 62, 91, 95, 99, 124, 131, 179]. A system for the
373 personalised computation of f_0 according to the sex and age of the speaker has been presented in [180],
374 using this value for the discrimination of normophonic and dysphonic voices. A different approach consists
375 on measuring the vocal function through the quantification of the energy contained in the signals. Since

376 this quantity is dependent on the distance between mouth and microphone, measures of SPL are preferred
377 instead. These can be achieved by means of intraoral or subglottal pressure apparatuses, or indirectly
378 computed by using accelerometers placed on the neck [77, 177]. Measurements of the vocal level at diverse
379 frequency ranges define voice range profiles (or phonetograms), which have been employed for voice pathology
380 detection [181, 182]. An extension to the method relies on the characterisation of a speech range profile,
381 that has also been used in AVCA systems [78, 183]. Finally, it is possible to extract features from the glottal
382 signal, characterising the opening and closing phases of the glottal waveform as in [52], or via the residue
383 obtained after inverse filtering, with measures such as the mean square residue or the excess coefficient
384 (kurtosis of the magnitude distribution) [69]. Other acoustic characteristics measure include, for instance,
385 the number and degree of voice breaks and unvoiced frames in speech [184].

386 4.3.2. *Perturbation and fluctuation analysis*

387 *Perturbations* are minor disturbances or temporary changes that deviate from an expected behaviour [22].
388 Perturbation parameters have been frequently used to analyse *vocal aperiodicities* that are the product of
389 modulation or additive noise. The most popular *modulation noise quantifiers* include families of parameters
390 based on jitter and shimmer.

391 On one hand, *jitter* is a measure of the short-term (cycle-to-cycle) perturbation of f_0 , with popular
392 examples including [185]: (i) *jitter relative*, which is the average absolute difference between consecutive
393 periods, divided by the average period; (ii) *jitter Relative Average Perturbation (RAP)*, which is the average
394 absolute difference between a period and the average of this and its two neighbours, divided by the average
395 period; and *jitter 5-point Period Perturbation Quotient (PPQ5)*, which is the average absolute difference
396 between a period and the average of this and its four closest neighbours, divided by the average period.
397 Jitter has been extensively used in AVCA systems, with some relevant examples in [52, 54, 66, 69, 85, 91,
398 95, 98, 99, 103, 131, 179, 184, 186–188]. On the other hand, *shimmer* measures short-term (cycle-to-cycle)
399 amplitude perturbations, having popular examples in [185]: (i) *shimmer absolute*, which is the variability of
400 the peak-to-peak amplitude in dB; (ii) *shimmer relative*, which is the average absolute difference between the
401 amplitudes of consecutive periods, divided by the average amplitude; and (iii) *shimmer 3-point Amplitude*
402 *Perturbation Quotient (APQ3)*, which is the average absolute difference between the amplitude of a period
403 and the average of the amplitudes of this and its neighbours, divided by the average amplitude. Literature
404 reports examples of shimmer in AVCA systems in [52, 85, 91, 95, 98, 99, 103, 131, 179, 184, 186–188]. Despite
405 both families of features have been of considerable utility for describing type I signals, its validity for type II
406 and III typologies has been put into question due to the need of an accurate identification of f_0 [189]. Some
407 methods have been devised to calculate jitter using spectral techniques, avoiding the need of a precise f_0
408 computation, having accomplished favourable results in pathology detection tasks [27, 148], whereas others
409 have derived shimmer and jitter through autoregressive decomposition and pole tracking [190]. Another
410 perturbation measure proposed for the characterisation of Parkinson’s disease is the *Pitch Period Entropy*
411 [107, 191, 192] which takes into account the smooth vibrato and microtremor in normophonic voices, and

412 the logarithmic nature of speech.

413 Popular *additive noise quantifiers* compute the relationship between the harmonics and background
414 noise contained in speech, with notable examples including the *Signal-to-Noise Ratio* (SNR) [82, 186–188],
415 *Harmonics-to-Noise Ratio* (HNR) [85, 91, 131, 187, 193] (and its variation Cepstral HNR [194, 195]), *Nor-*
416 *malized Noise Energy* (NNE) [42, 95, 187, 196] (and its variation adaptive-NNE [135]), and *Glottal-to-Noise*
417 *Excitation Ratio* (GNE) [138, 145, 194]. These features have been extensively applied in the evaluation of
418 voice quality correlating positively with many speech disorders and perceptual ratings [37, 138]. Other noise
419 quantifiers include the empirical mode decomposition excitations ratios, the vocal fold excitation ratios or
420 the glottis quotient, which have been employed for the characterisation of Parkinson’s disease [191, 192].

421 By contrast, *fluctuation analysis* is referred to the study of severe disturbances in the dynamics of
422 the vocal folds behaviour, reflecting the inherent instability of the system [22]. *Tremor* is one prominent
423 characteristic that is often studied, referred to low-frequency fluctuation in amplitude and/or frequency
424 (modulation noise), related to pathologies of neurologic origin [22]. Some popular tremor estimates include
425 *amplitude and frequency tremor* [179, 197, 198], or the *turbulent noise index* [58, 122].

426 4.3.3. Spectral-Cepstral analysis

427 Measures derived from the acoustic spectrum/cepstrum have been widely used in the study of pathologi-
428 cal phonation, to mostly characterise *vocal quality*. Indeed, spectral and cepstral features have demonstrated
429 high correlation with the perceptual assessment of dysphonia, providing large sensitivity when used in classi-
430 fication tasks [11]. The analysis with spectral/cepstral features also presents several advantageous properties,
431 including its appropriateness for analysing both sustained vowels and running speech with no extra proce-
432 dures, and the ability to characterise speech signals without depending of the estimation of f_0 [11, 37].

433 Spectral measures derived directly from the speech spectrum include diverse *Long-Term Average Spectrum*
434 (LTAS) characteristics [23, 28, 76, 85], where the spectral tilt [23] -which indicates the degree of energy
435 concentrated in low- vs. high-frequency areas of the spectrum- is a noteworthy example. Other popular sets
436 of features are based on the estimation of the spectral envelope with a noteworthy example in the *Linear*
437 *Predictive Coding Coefficients* (LPC) [126, 153, 199–202], which can also be used to decompose the speech
438 signal into its residual and vocal tract components and hence derive parameters for characterisation [153].
439 Other features extracted from the residual signal include the *pitch amplitude* [23, 28, 42], which measures
440 the dominant peak in the auto-correlation function of the residual signal; and the *spectral flatness ratio*
441 [28, 42, 69], which measures the residual flatness. Some variations of LPC include cepstral transformation like
442 the *Cepstral Linear Predictive Coding Coefficients* (LPCC) [48, 121, 141, 201–203], or a mel-transformation
443 of LPC called Mel-line spectral frequencies [169].

444 Literature also reports several approaches relying on filter-banks to decompose voice signals into different
445 sub-bands. For instance in [73], correlation functions measure the relationship between the bands of an octave
446 filter-bank for the detection of pathologies. Another example is in [117], where HNR is computed on different
447 frequency bands along with the energy. It is also possible to employ filter-banks that rely on psychoacoustic

448 criteria to condense spectral information into perceptually-relevant frequency bands. Within this category,
449 *Mel-Frequency Cepstral Coefficients* (MFCC) are probably the most popular features in speech-related ap-
450 plications, exploiting auditory principles and the decorrelating property of the cepstrum to characterize
451 speech signals. MFCC have been used extensively in several works [43, 57, 115, 121, 130, 133, 140, 141,
452 146, 152, 156, 169, 201, 203–205], and often in combination with its first derivative (Δ) representing velocity
453 [81, 150, 206], and/or its second derivative ($\Delta\Delta$) representing acceleration [116, 128, 129, 139, 147, 163].
454 Works related to MFCC characterisation include the study of the influence of the tapering window that
455 is employed for the spectrum estimation during the MFCC processing [114], or the derivation of an in-
456 dex based on MFCC: the *pathological likelihood index* [157]. Other psychoacoustic characteristics include
457 *Perceptual Linear Predictive Coding* (PLP), which have been designed in accordance to a scale modelling
458 the human auditory system [207–209]. This has been used in AVCA systems in [48, 163, 210], next to
459 its bandpass filtered variation RASTA-PLP [48, 141, 163, 201, 209, 210]. Other characteristics include the
460 band power decorrelation time obtained through the Meddis and O’Mard filterbank model [37]. The filter-
461 banks can also be employed to perform time-frequency decompositions, through wavelets transformations
462 [50, 57, 60, 122, 123, 134, 165, 199, 211–215] -which are often accompanied by the calculation of energy
463 and/or entropies of the sub-band decompositions- or adaptive time-frequency distribution [164]. An inter-
464 esting review about these time-frequency decompositions is presented in [216].

465 Finally, measures based on the cepstral prominence of the harmonics are often utilised, by means of
466 features such as *Cepstral peak prominence* (CPP) (next to its variation called smoothed CPP), which is
467 a normalised measure of the cepstral peak amplitude, comparing the level of harmonic organisation of the
468 speech to the cepstral background noise resulting from aspiration [217, 218]. This measure has been reported
469 as one of the strongest correlates of breathiness [11, 15, 37]. Several studies have indicated that cepstral
470 measures may be supplemented by other acoustic quantifiers such as the *Low-to-High Harmonic Ratio* (LHr)
471 which measures the spectral tilt of the spectrum above and below 4 kHz [15]. Relevant examples employing
472 these measures include [11, 25, 76, 80, 85, 97, 186]. A derived index that incorporates the information
473 provided by CPP, LHr, and standard deviation of LHr is the cepstral spectral index of dysphonia presented
474 in [75].

475 4.3.4. Complexity analysis

476 *Complexity* is a controversial term that has been classically linked to randomness and mistakenly asso-
477 ciated to information measures such the algorithmic complexity [219]. However, it is more appropriately
478 related to a "*meaningful structural richness*" [220] or to *fractal* behaviour rather than to randomness. Com-
479 plex behaviour is typically observed in biological systems that manifest at least one of the following dynamical
480 properties [221]: (i) nonlinearity; (ii) nonstationarity; (iii) time irreversibility (or asymmetry); or (iv) mul-
481 tiscale variability. One of the most popular approaches to investigate the complexity of a system is through
482 *Nonlinear Dynamics Analysis* (NDA). Nonlinear phenomena arises naturally in physiological systems, and
483 voice production is not an exception to this regard. Indeed, supporting findings of nonlinearity in phonation

484 include nonlinear pressure-flow relations in the glottis, the delayed feedback from the mucosal wave, the
485 nonlinear stress-strain curves of vocal fold tissues, nonlinearities associated with vocal fold collision [222],
486 or asymmetries between the right and left vocal folds vibrations [223]. In addition, nonlinear dynamical
487 behaviour of models of the vocal folds such as period-doubling (subharmonics), bifurcations, and transitions
488 to irregular vibration have been observed in experiments with excised larynges; whereas period-doubling
489 bifurcations and chaos-like features have been observed in signals from patients with organic and functional
490 dysphonia [224]. Aforementioned facts suggests the appropriateness of using NDA to characterise the dy-
491 namics of voice production even in pathological scenarios, as voice pathologies can be considered disorders
492 of glottal dynamics [225]. Indeed, the complexity features attempt to measure *vocal aperiodicity* phenomena
493 in AVCA systems. The usual approximation to NDA relies on a reconstruction, termed *embedding*, to reveal
494 the system dynamics in a space called *state space*. The most popular indices calculated using NDA, compute
495 the dimensionality of the reconstructed state space. They have been used for discrimination of pathological
496 and normophonic states, being popular examples the fractal dimension [72] or the *correlation dimension*
497 (D2) [56, 67, 98, 103, 107, 156, 167, 186, 188, 226]. Other measures are based on the rate of divergence of
498 trajectories in the state space. This has been explored to differentiate normophonic and dysphonic voices
499 through the computation of *Largest Lyapunov Exponent* (LLE) [54, 71, 94, 141, 156, 167, 226, 227] or the
500 Lyapunov spectrum [166]. Measures of entropy (the rate of information gain) have also been employed in
501 AVCA systems, by means of the first and/or second order entropies [68, 98, 188], the relative entropy [61],
502 pseudo-estimators such as the Ziv-Lempel complexity [71, 141, 226], or entropy metrics based on *Hidden*
503 *Markov Models* (HMM) [156, 228]. A concept related to entropy is the *regularity*, which measures the pre-
504 dictability of the time series. The most popular regularity estimator used in pathological voice analysis is
505 the ApEn [87, 88, 149, 226]. Other ApEn-derived metrics used in AVCA systems include the SampEn [51],
506 GSampEn and FuzzyEn [71, 156, 228].

507 Likewise, time-frequency decomposition have been employed to explore the fractal properties of speech
508 [176, 229] or to characterise complexity (using ApEn) on each decomposed sub-band [212]. Other measures
509 explore the use of nonlinear prediction [70, 227]; measure the self-similarity of the voice by means of the
510 DFA [107, 141, 191, 192, 224], the Hurst exponent [71, 94, 141, 176, 212]; or characterise properties of
511 recurrence to compute the effects of modulation noise using the *Recurrence Period Density Entropy* (RPDE)
512 [107, 141, 191, 192, 224, 230]. A discussion about characterisation using nonlinear dynamics can be found
513 in [231].

514 4.3.5. 3-dimensional representations

515 A popular approach that has been gathering attention lately is based on the multidimensional represen-
516 tations of speech, and the employment of image processing techniques or matrix tools for the extraction of
517 characteristics. In this regard, one approach is based on *Modulation Spectra* (MS), which characterises the
518 modulation and frequency components of speech. MS produces a 3-dimensional representation that has been
519 employed for the detection of voice pathologies [119, 140, 232].

520 Following the idea of 3-dimensional representations, a mel-spectrogram is characterised by means of
521 the recurrence texture plots and the local binary pattern operator in [233]. Matrices can also be formed
522 using time-frequency decomposition as in [132], where features are extracted through the employment of a
523 multidirectional regression, the use of interlaced derivative patterns of the glottal source excitation signal
524 as in [120]. Similarly, authors in [160] extract features from co-occurrence matrices formed after using
525 filter-banks on the input speech or octave-spectrogram [142].

526 4.3.6. Other types of features

527 There exist some other approaches that do not fit into the above categorisation. For instance, the multi-
528 dimensional *acoustic voice quality index* is a metric based on weighted multivariate regression of 6 temporal,
529 spectral and cepstral characteristics that has been used for voice pathology assessment and detection [234].
530 The use of variograms and the characterisation with the signal-to-dysperiodicity ratio has been explored in
531 [25, 116]. Likewise, decompositions of speech based on non-negative matrix factorisation [164] or empirical
532 mode decomposition [168] have also been reported. The hoarseness diagram has been proposed to visu-
533 alise additive and modulation noise components [53, 145, 235]. The utilisation of higher order statistics for
534 the characterisation of dysphonic voices is discussed in [137], whereas the spectral properties of centralized
535 auto-correntropy is used to detect and classify vocal pathologies in [236].

536 Some multidimensional approaches consider a combination of several features of different type, having
537 some notable examples in [28, 55, 65, 96, 144, 156, 171, 172, 237, 238]. Other multidimensional approaches
538 consider metrics derived from the MPEG-7 standard as described in [128, 154, 162]. Another type of character-
539 isations are based on biomechanical models to describe the behaviour of the glottal and mucosal waveforms
540 [172, 237]. Finally, literature also reports the employment of MFCC features for the construction of a
541 phonological model of 14 features (voicing, place of articulation, turbulence, nasality, etc.) [86].

542 4.4. Dimensionality reduction

543 Dimensionality reduction is aimed to decrease the size of the feature space in order to remove redundant
544 or irrelevant features that might affect performance. Two major types of techniques can be defined: those
545 based on feature extraction, which employ a transformation of the input space; or those based on feature
546 selection, not relying on any transformation. *Feature extraction* techniques include the classical approaches of
547 singular value decomposition [134], *Linear Discriminant Analysis* (LDA) [123, 161, 166, 199, 211, 237, 239],
548 and *Principal Components Analysis* (PCA) [25, 57, 123, 124, 127, 136, 145, 153, 161, 172, 227, 237–239], or
549 extensions of PCA such as kernel-PCA [127, 238], neural-networks PCA [127], or dynamic feature extraction
550 using PCA [150]. Other type of transformations include those based on HMM [71, 136], clustering-based
551 feature weighting methods [134, 161] or others based on multiple regression analysis [93]

552 Within *feature selection*, two types of methods arise. In one hand, the *wrapper feature selection* ties
553 the selection of features to the maximisation of a performance metric obtained with a classifier/regressor.
554 Some notable examples in AVCA include the use of genetic algorithms to select the best set of features for

555 recognition purposes [50, 57, 90, 166, 212], binary logistic regression analyses with stepwise variable selection
556 [78], sequential backward and forward feature selection [65, 239] or angle modulated differential evolution
557 [63]. On the other hand, *filter feature selection* employs correlation and information approaches to find
558 the most pertaining sets of features. In this respect, literature reports the use of the mutual information
559 [145, 171], correlation analysis [145], Fisher discrimination ratio [47, 129, 129, 154, 226, 233, 237], Fisher
560 discriminant analysis [144], or the Davies-Bouldin index [212].

561 4.5. Machine learning and decision making

562 The machine learning procedure receives different names depending on the type of decision making task
563 that is involved. If given a set of observations $\mathcal{X} = \{\vec{x}_1, \dots, \vec{x}_n, \dots, \vec{x}_N\}$, where each \vec{x}_n is associated
564 to a label $\vec{\ell} = \{\ell_1, \dots, \ell_n, \dots, \ell_N\}$, the aim of the procedure is to learn a mapping from the input set of
565 observations to the labels. This type of task is known as *supervised learning*; in opposition to *unsupervised*
566 *learning*, which is related to the discovery of structure in the data in the absence of labels.

567 To the author’s knowledge all of the machine learning methodologies presented in AVCA systems are
568 based on supervised learning. Within this category, the most widely employed decision machines include the
569 *Support Vector Machines* (SVM) or *Gaussian Mixture Models* (GMM). SVM is a discriminative classifier
570 constructed from sums of kernel functions which has been used in AVCA systems in [47, 50, 52, 60, 62,
571 63, 71, 73, 86, 90, 91, 94, 96, 116, 119, 120, 123, 140, 141, 150, 153, 154, 166, 169, 199, 201, 210, 212,
572 226, 229, 233, 239–241]. By contrast, GMM is a type of generative classifier that has provided excellent
573 results in diverse speech-related applications. Its popularity arises from the modelling capability they offer
574 and the probability framework on which they stand. The use of GMM in AVCA systems is reported in
575 [114, 121, 126, 129, 132, 139, 140, 150, 156, 169, 200, 201, 204], next to variations such as GMM-Universal
576 Background models [81, 163], GMM-SVM [130, 133, 201], *i*-Vectors [163].

577 Other popular decision machines include *Artificial Neural Networks* (ANN) [49, 57, 68, 123, 137, 143,
578 146, 152, 165, 166, 186, 211, 239], *Deep Neural Networks* (DNN) [142, 151, 155], HMM [74, 115], random
579 forests [63, 65, 65, 94, 187], LDA [23, 25, 42, 122, 131, 139, 168, 169, 179, 203, 213, 214] quadratic-LDA [224],
580 Hidden Markov Models [52, 71, 115, 121, 136, 139, 201, 228], *k*-nearest neighbours (KNN) [96, 117, 203, 212]
581 and the Bayes classifier [150]. The use of regression techniques has also been reported in [28, 75, 78, 85, 234].

582 4.6. Evaluation of the system

583 A common approach followed in machine learning applications to generalise results and provide valid
584 measures about the actual efficiency of the systems consists in the use of validation techniques. The basis of
585 these methodologies is the decomposition of the available dataset into subsets which are used independently
586 for training and testing purposes (and often for parameter tuning). On one hand, the *training partitions* are
587 used to estimate a mapping between observations and labels in the supervised machine learning algorithms.
588 On the other, the *testing partitions* are employed to assess the performance of the system. A third partition
589 often arise for the purposes of parameter tuning of the algorithms.

590 The most straightforward approach for validation is the *split sample* method, which consist on using a
591 percentage of the dataset to conform exclusive training and testing partitions. The problem of split sample
592 is the lack of generalisation of results -specially if the data is scarce- as the randomly sampled partition
593 might not be representative of the data under analysis, and the reported results will probably depend on the
594 partitions that have been chosen. The use of this methodology has been reported in [52, 123, 137, 147, 169,
595 184, 201, 206, 239].

596 A different approach -which is one of the most popular validation techniques- is the *k-folds cross-*
597 *validation*, which produces k disjoint sets of size N/k , called folds, with N representing the number of
598 observations. In total k iterations are performed, using in each case a different subset for testing pur-
599 poses, and the remaining $k - 1$ for training algorithms. The measures of performance are then evaluated
600 as the mean value calculated among iterations. Works using cross-validation techniques include [43, 47-
601 49, 55, 57, 63, 63, 68, 71, 73, 91, 119-121, 124, 126, 127, 130-133, 136, 140, 146, 150-156, 160, 161, 168, 199,
602 200, 204, 211, 229, 233, 236, 240-242].

603 Another popular evaluation methodology is the *leave-one-out validation* which arises in the limit $k = N$
604 in a *k-folds cross-validation*. In this case, only one observation is used for testing and the remaining registers
605 are employed for training, repeating this procedure N times. Leave-one-out validation is usually preferred to
606 cross-validation when the dataset sizes are small as it allows to maximise the size of the training partition.
607 Notable examples of leave-one-out validation are reported in [23, 49, 81, 96, 176, 213, 214].

608 In the same way, *bootstrapping* consists on randomly selecting a number of points from the training
609 partition, with replacement, to train machine learning models and then calculating performance on a testing
610 partition. This process is repeated k times, thus generating k different models. At the end, the performance
611 is computed as the mean performance obtained in the testing partition. Some examples reporting the use of
612 bootstrapping in AVCA are presented in [65, 75, 89, 101, 136, 224].

613 In addition to those validation methodologies which are trained and tested in the same corpus, other
614 manners to validate performance are based on *cross-dataset validation* on which the training set corresponds
615 to a particular dataset, whereas the testing partition corresponds to another different one. The advantage
616 of such an approach is in the possibility of testing the robustness of the AVCA system in a more realistic
617 scenario with an increased variability. This methodology has been used in works such as those presented
618 in [47, 120, 148, 151, 204, 243].

619 With regards to metrics of performance, the simplest approach consists on computing measures that
620 compare the predicted labels given by the decision machines to the actual labels of the dataset. In this
621 regard the most commonly used metric is the *accuracy* (ACC) -which has been used in almost all the
622 reviewed papers- representing the rate of the correctly identified labels in comparison to the total number of
623 instances. Another manner to analyse the performance of binary detection systems, is by means of *Receiver-*
624 *Operating Curve* (ROC) [244] and *Detection Error Tradeoff* (DET) curves [245]. An additional measure
625 derived from these curves is the *Area Under ROC Curve* (AUC), which is a value ranging between 0 and 1,

626 that is obtained after integrating the ROC curve. A number of reasons favour the use of this metric instead
627 of other classical measures such as the ACC, including [246]: (i) a standard error that decreases as both
628 AUC and the number of test samples increase; (ii) decision threshold independence; (iii) and invariance to
629 a-priori class probabilities. Some examples demonstrating the use of ROC curves and AUC are included in
630 [23, 129, 144, 164, 166, 201]

631 Other types of performance evaluation techniques include the cost of log-likelihood-ratio [63, 89, 158], the
632 DET curves [43, 65, 89, 91, 119, 129, 133, 138, 140, 146, 150, 150, 152, 174], or sensitivity versus specificity
633 plots [75]. Statistical analysis based on the Mann-whitney U-test [69, 98, 103, 126, 167, 179, 188] or the t-test
634 [47, 51, 57, 59, 64, 73, 80, 146, 187] have also been reported to compare the means of different populations.

635 4.7. Applications of AVCA systems

636 There exist a variety of applications of AVCA methodologies to characterise a wide variety of voice
637 impairments. Without being extensive, the following will introduce some relevant applications of AVCA
638 systems. In this manner, most of the works in literature are related to the analysis of laryngeal pathologies
639 such as nodules [97], polyps [59, 98], larynx cancer [87, 88, 99, 149], diplophonia [101, 178], spasmodic
640 disphonia [102], unilateral laryngeal paralysis [50, 103], laryngectomised patients using oesophageal voice
641 [70]. Notwithstanding, there exist several works focusing on other disorders such obstructive sleep apnea
642 [104, 105, 247–249], hypernasality [106, 141, 250], Parkinson’s disease [92, 107, 163, 192, 209, 210], Alzheimer’s
643 disease [251, 252], dysphagia [109, 110], lupus [111].

644 5. Aspects affecting AVCA systems

645 The variability embedded in speech has long been recognised as a major source of errors in automatic
646 classification systems based on speech. For instance, several variability factors identified in the design of
647 speaker recognition systems are described in [253]. Translated to terms of AVCA systems, these variability
648 factors are the following: (i) *Peculiar intra-class variability*: manner of speaking, age, sex, inter-session
649 variability, dialectal variations, emotional condition, etc. (ii) *Forced intra-class variability*: Lombard effect,
650 external-influenced stress, cocktail-party effect, etc. (iii) *Channel-dependent external influences*: type of mi-
651 crophone, bandwidth and dynamic range reduction, electrical and acoustical noise, reverberation, distortion,
652 etc. It is worth noting that the forced intra-class variability is more common of unsupervised recording
653 environments rather than from controlled clinical settings. That does not imply, though, that their effects
654 should be disregarded. For instance, they are of significant importance in telemedicine scenarios where the
655 recording conditions might vary widely. Despite that, and for the purpose of simplicity, the forced intra-class
656 variability is to be omitted from further discussions, and the term *intra-class variability is to be referred*
657 *to the peculiar intra-class variability only*. Moreover, the intra-class and channel-dependent factors can be
658 further associated to the linguistic, paralinguistic, transmittal and extralinguistic spheres.

659 Having this in mind, the current section introduces some factors that are of interest in the design of AVCA
660 systems. It is worth noting that this list is not exhaustive but it only include what could be considered as
661 the most important factors affecting these systems.

662 5.1. Intra-class variability

663 One important aspect to be outlined in any speech-related system is the effect of the intra-class or
664 *intersession* variability. In a speaker recognition system, aiming at recognising the identity of a target
665 speaker, the intersession effect might be described by the differences arising between recordings of the target
666 speaker due to vocal effort, physical or emotional condition, etc. In the AVCA field, the intersession variability
667 might be explained by the acoustic diversity among different pathologies, the sex or age of the speaker, or the
668 spurious information introduced by other linguistic, paralinguistic or extra-linguistic effects. In this regard,
669 several intra-class factors that might affect the performance of voice pathology classification systems include
670 *linguistic* aspects such as the speech production task, the dialect and accent of the speaker; *paralinguistic*
671 events such as the emotion or the vocal effort; or *extralinguistic* effects such as the sex or age of the speaker.

672 5.1.1. Dialects and accents

673 *Dialects* are the result of systematic, internal linguistic changes that occur within a language, reflected in
674 the form of structural alterations in phonology, morphology, syntax, lexicon or semantic [33]. Dialects have
675 been identified as an important aspect when defining communication disorders. Undoubtedly, not accounting
676 for dialect features may result in the misdiagnosis of communication disorders [33]. For instance, several
677 key features of African-American English phonology have been found to overlap with identifiers of speech
678 sound delay/disorder in the phonology of general American English, making the distinction from normal and
679 disordered states problematic in African-American speakers [254]. This phenomenon has been found in other
680 contexts where non-prestige social dialects are often incorrectly associated to disordered speech [33]. *Accents*,
681 on the other hand, are linguistic changes within a language that occur mainly at the phonological level. It has
682 been long identified as an important confounding source in speech-related applications. For instance, accent
683 is described as the most important source of variability between speakers in speech recognition systems in
684 [255]. A further study presented in [256], confirmed that accent degrades classification rates, with errors
685 increasing around 40-50% in cross-accent speech recognition scenarios. In general, it has been found that
686 performance degrades when recognizing accented from non-native speech [257].

687 5.1.2. Vocal effort

688 *Vocal effort* is a subjective physiological interpretation of the voice level, as given by judges, or by the
689 speaker itself to adapt speech to the demands of communication [258]. There exist evidences indicating that
690 the vocal effort alter perceptual and acoustic parameters extracted from speech, and therefore might impact
691 AVCA systems. Phenomenologically, variations in vocal effort affect the shape of the glottal pulses, changing
692 the closing velocity waveform and affecting the relative duration of the closed interval [258]. In voices

693 produced with increased vocal effort there have also been found significantly greater values in parameters
694 such as subglottal pressure, translaryngeal airflow, and maximum flow declination [259]. Similarly, it has
695 been stated that the medial compression of the vocal folds is enhanced when vocal effort is augmented, which
696 results in an improved glottal closure, enlarged vocal intensity, and increased f_0 and amplitude [260]. Vocal
697 effort also alters the duration of vowels, consonants, and the pausing behaviour during speech production
698 [258]. In terms of quality, voices produced at excessive vocal effort are perceptually described as creaky
699 [258] or strained [259]. Not surprisingly, this has consequences on parameters extracted from speech. For
700 instance, jitter, shimmer, NNE and two EGG parameters have been found to vary significantly among vowels
701 produced at three vocal effort levels (low, normal, high) [261].

702 Similarly, it has been reported that jitter and shimmer significantly increased their value with decreasing
703 voice intensity [262], being also identified as one of the most important factors influencing the computation
704 of these perturbation parameters, alongside with the sex of the speaker and the type of uttered vowel [263].
705 Other sets of parameters which are affected include cepstral features, which have been reported to differ
706 substantially at diverse effort levels [260]. This has been confirmed in [259], where significant differences
707 arose in 4 aerodynamic and 2 cepstral measures when comparing phonation at different effort levels. Authors
708 in [264] have investigated the effects of increased vocal effort in pathological phonation. Results indicate
709 that louder voicing reduces the values of perturbation parameters in normophonic speakers or in superficial
710 vocal fold pathologies, while in cancer or vocal fold paralysis, louder phonation significantly enhances the
711 irregularity of vocal folds vibration.

712 5.1.3. *Emotion*

713 The study of the emotional content embodied on speech has garnered a lot of attention within the speech
714 research community. Indeed, the term *affective computing* has being coined to describe the automatic
715 sensing, recognition and synthesis of human emotions from any biological modality such as speech or facial
716 expressions [265]. There exist some studies considering the effect of emotions in speaker recognition systems.
717 For instance, emotions are regarded as a factor affecting automatic recognition of children's speech in [266].
718 In [267], the effect of an emotions recogniser previous to a speech classification process is investigated.
719 Authors report that affective speech downgrades the system performance, and that a cascading scheme is
720 highly effective in improving recognition rates. Despite these facts, little is known about the influence of
721 emotion in AVCA schemes, but according to the evidence found in the field of speaker recognition, it might
722 be hypothesized that affective speech is a confounding factor that should be taken into consideration.

723 5.1.4. *Sex*

724 The variability introduced by the sex of the speaker remains as a major concern in the design of speech-
725 based systems. Indeed, authors in [255] reported that this factor accounted for the second most prominent
726 source of variability -after accent- in speech recognition systems. Certainly, literature states that the perfor-
727 mance of speech recognition, identification or verification systems improves by employing a-priori information

728 about the sex of the speaker [268]. For instance, authors in [269] obtained a 2% of accuracy improvement in
729 a speaker recognition system when using sex-specific models.

730 The nature of the variability introduced by the sex of the speaker stands on physiological, acoustic, and
731 psychophysical factors [6]. Regarding *physiological* differences, the human laryngeal anatomy differs between
732 sexes at a variety of levels. Particularly, males tend to have a more acute thyroid angle; thicker vocal folds; a
733 longer vocal tract; a larger pharyngeal-oral apparatus, thyroid lamina and skull compared to that of females
734 [270, 271]. Studies of excised human larynges have shown that anteroposterior dimensions of the glottis are
735 1.5 times larger in men than in women [272]. Besides that, the female pharynx has been found to be shorter
736 than of males during the production of the three cardinal vowels. This may be a key factor in distinguishing
737 between male and female voice qualities during speech production [271]. In addition, the observation of the
738 glottis during phonation has suggested the presence of a posterior glottal opening that persists throughout
739 a vibratory cycle and which is common for female speakers, but occurs much less frequently among male
740 speakers [273]. Indeed, about 80% of females and 20% of males have a visible posterior glottal aperture during
741 the closed portion of a vocal period [274]. Regarding *perceptual* differences, parameters such as effort, pitch,
742 stress, nasality, melodic patterns of intonation and coarticulation are used for characterising female voices,
743 while male voices are judged on the basis of effort, pitch and hoarseness [275]. It is also argued that female
744 voices possess a "breathier" quality than male voices [274]. The pitch is the most known trait differentiating
745 sexes [275], with females' pitch higher than of males' by as much as an octave [276]. This pitch difference
746 might influence the perception of dysphonic voices since lower pitch is perceived as rougher [31]. In addition
747 to the pitch, literature reports significant differences between male and female speakers' formants (f_1 , f_2 ,
748 f_3 , f_4) [6]. This is because the vocal tract length for males is longer than that of females, producing on
749 average formant patterns scaled upward in frequency by about 20% [275]. There are also several important
750 acoustic consequences of the posterior glottal opening during the closed phase of phonation, which is more
751 frequent in females. A first consequence is a breathier voice quality which is the result of a larger amount
752 of air passing through the vocal tract [270] and that affects the relative amplitude of the first harmonic of
753 the speech spectrum [272, 277]. A second consequence is the widening of the f_1 bandwidth, which is the
754 result of the glottal aperture that produces energy losses particularly at low frequencies [273, 277]. A third
755 acoustic consequence is the generation of turbulence in the vicinity of the glottis [277], perceived with a high
756 level of aspiration noise in the spectral regions corresponding to f_3 , and contributing to a breathier voice
757 quality [276]. A final consequence is a lower spectral tilt due to the presence of aspiration noise [276], which
758 turns out to be a significant parameter for differentiating between male and female speech samples [273].

759 In addition to the acoustic differences reported from the study of the raw speech, there are some differences
760 in the glottal components among sexes. On one hand, the female glottal waveform tends to have a shorter
761 period, lower peaks and peak-to-peak flow amplitudes than that of males [278]. Likewise, the derivative of
762 the glottal waveform does not present an abrupt discontinuity during the closing time due to the incomplete
763 closure of the vocal folds [6]. In general, it is stated that female glottal components are symmetric, with

764 opening and closing portions of the waveform tending towards equal duration [279]. Conversely, and regarding
765 the glottal waveform of male speakers, it is found that the open quotient is smaller and the maximum flow
766 declination rate is greater than of females [273]. Moreover, the closing portion of the waveform generally
767 occupies 20 – 40% of the total period and it might not exist an easily identifiable closed period [275]. In
768 general, it is stated that male glottal waveforms are asymmetrical and present a hump in the opening phase.
769 Finally, it is worth noting that sex differences are found to be age and hormone-dependent, and thus is of
770 great importance considering the effect of age when studying male or female voices.

771 5.1.5. Age

772 According to the male-female coalescence model of ageing voice, hormone-related factors cause changes
773 in voice production systems. In this manner, hormones during puberty are responsible for the differences
774 between males and females in adolescence, but these changes are counteracted to some degree by hormone
775 related factors associated with menopause and ageing [271]. During males' puberty, the thyroid cartilage
776 develops the Adam's apple, the muscular and mucosal layers of the vocal folds thicken, the vocal folds lengthen
777 and widen, the cricothyroid membrane widens, and the corresponding muscle becomes more powerful [270,
778 280]. As a result, the fundamental frequency decreases an octave compared to that of a child [280]. During
779 females' puberty, there is little development of the thyroid cartilage or of the cricothyroid membrane, and
780 the vocal muscle thickens slightly but remains supple and narrow. As a result, the female's f_0 becomes one
781 third lower than that of a child [280]. The age effects on the larynx tend to be more significant in men than
782 in women. In this manner, males experience an increasing of the fundamental frequency as a result of muscle
783 atrophy, thinning of the lamina propria, general loss of mass and ossification and calcification of the larynx
784 that starts during the third decade of life [271]. In females, ossification and calcification starts in the fourth
785 decade of life, and in some cases never completely ossify. However, due to menopause effects, a lowering of
786 fundamental frequency prior to senescence occurs [271].

787 The effects of age in AVCA systems have not been studied in depth. One of the few works that accounts
788 for its effects in AVCA is introduced in [75], where this trait has been used as a predictor in a binary logistic
789 regressor for the prediction of dysphonia, having found a marginal but statistical significant increment in
790 performance when age is introduced in the model.

791 5.2. Channel-dependent external influences

792 This dimension includes all the effects that aggregate variability to speech registers because of the mere
793 act of recording. This is a well-known problem that has long been identified in speech and speaker recognition
794 systems [113]. Several aspects affect the recording process, including the instrumentation (type of micro-
795 phone, analogue-to-digital converter, etc.), the acoustic environment (office, recording studio, etc.) and the
796 transmission means (land-line, cellular, etc.) [281]. Similarly, background noises, noises made by speakers
797 (such as lip smacks), noise in the input device itself, etc., are recognised as sources that impair performance
798 of speech recognition systems [113]. Another problem that might arise, is the variability introduced because

799 of the mismatch in the recording conditions, between the registers employed for training the models and
800 those used for testing purposes. For instance, as when a certain microphone is used for recording the train-
801 ing utterances, but the model is verified with a different equipment. Indeed, the microphone is expected to
802 modify the speech spectrum, and anything that modifies the spectrum may cause difficulties in recognition
803 tasks [281]. In this regard, the study of [38] demonstrated that the type of microphone has effects in the
804 computation of perturbation parameters, providing evidence favouring the use of condenser cardioid micro-
805 phones instead of dynamic or omnidirectional microphones. This study also showed that sensitivity and
806 microphone-to-mouth distance have the largest effect on perturbation measures, whereas the angle had little
807 effect for short distances, but a greater effect for longer distances. The study presented in [282] showed that
808 a signal-to-noise ratio of 42 dB is needed to provide reliable estimations of perturbations measures, whereas
809 values less than 30 dB have been shown to impact negatively in their computation. In [283], the effect of
810 background noise, reverberation, clipping and speech compression on the calculation of MFCC features was
811 tested out, demonstrating significant (but predictable effects) in MFCC computations.

812 6. Discussions and conclusions

813 This paper has presented concepts in relation to voice impairments and AVCA systems. In this manner, a
814 categorisation of diverse aspects of voice conditions in terms of perceptual and physiological phenomena has
815 been proposed, as well as a description of a prototypical AVCA system along with each one of its constituting
816 parts. With relation to the latter, a systematic review has been carried out to overview the methodologies
817 that are more often employed in AVCA systems.

818 Regarding the categorisation presented in section 3 and according to the systematic review of section
819 4, some inferences can be made. Firstly, a large number of papers still employ the MEEI corpus despite its
820 well-known limitations. Notwithstanding, the field has received the advent of novel public datasets such as
821 the SVD or privative corpora shared among different research groups, which has permitted the reproducibility
822 of results and has opened up the possibility to carry out comparisons among methodologies in other datasets
823 apart from MEEI. Despite that, there is room for improvement, as there is still necessary to record larger
824 datasets, more balanced in terms of pathologies, age or sex, and containing a larger variety of acoustic
825 material based either in sustained vowels, isolated words or running speech.

826 Literature has also revealed that most of the works employ sustained phonation due to its simplicity,
827 despite the potential that running speech presents. In this regard more investigation on novel features and
828 methodologies that employ this type of speech task is required. Similarly, it has been found that the effects
829 of extralinguistics or paralinguistics (such as age, sex, accent, etc.) have been seldom considered in the vast
830 majority of systems reported in literature, despite their relevance as confounding factors on this type of
831 systems. Accounting for this variability factors should be a relevant matter to study in the future.

832 Regarding the characterisation techniques, most of the reviewed papers report the employment of de-
833 scriptors to quantify vocal quality. However, vocal aspects describing variations in intensity and f_0 are also

834 important aspects to consider as they might serve to characterise other phenomena such as hypophonia or
835 inadequate pitches in speakers. Quantifying both intensity and f_0 might complement the information ob-
836 tained with descriptors of voice quality, with potential -for instance- to improve results of differential analysis
837 in identification tasks. This comes with the added cost, though, of having to record other variables besides
838 audio, such as the SPL or EGG (or to employ robust f_0 estimators such as those based on inverse filtering).
839 In the same manner, the systematic review also indicates that the increasing need of novel biomarkers for the
840 early diagnosis, differential analysis or assessment of disorders such as Parkinson, Alzheimer or obstructive
841 sleep apnea, has generated an emerging interest on quantifying dysphonic conditions. The analysis of these
842 disorders have also brought new features and processing techniques which have indeed enriched the field.

843 Regarding the machine learning and decision making methodologies, there is a large amount of method-
844 ologies related to supervised learning. The typical approach followed in AVCA systems is based on bottom-up
845 schemes on which the voice disorders phenomena is firstly studied, to build up systems from the inferences
846 obtained in the previous analysis stage. Notwithstanding, other related fields (such as speech recognition)
847 have experienced an increased interest in up-bottom schemes through unsupervised methodologies for the
848 purposes of pattern discovery or data mining. This same path could be followed in AVCA systems as well.

849 The literature review also served to reveal the existence of certain methodological issues that might
850 compromise the interpretation and validity of results in certain works. In an attempt to provide general
851 recommendations, some practical considerations are indicated next. One concern that is found in a few
852 papers is due to the addition of registers of other corpora besides the one under study, for the purposes of
853 balancing patients in terms of pathology, sex, age, etc.; or simply to increase the size of the studied corpus.
854 The effects of following such an approach are certainly important and might bias or affect the validity of
855 results due to the channel divergences between the datasets. Likewise, and despite there is an trend on
856 employing more robust validation techniques, there are certain issues that should be taken into account as
857 well. For instance, there are some works that suggest having included registers of the same speaker on both
858 training and testing partitions. This seems to be the results of having datasets on which speakers record
859 audio in different sessions, but not accounting for the possibility of including audio of the same speaker in
860 a training or a testing partition. Following this approach might introduce speaker information that might
861 bias the machine learning algorithms. Besides that, one problem that is common in a variety of papers is
862 to report results with confidence intervals which are larger than the range of measurement. For instance,
863 there are certain works reporting accuracy values of the type $99 \pm 1.5\%$. It is recommended the use more
864 robust estimators of confidence not to allow values larger than the range of acceptable values. A few papers
865 have reported the employment of energy measurements to characterise normophonic or dysphonic conditions
866 without having used SPL or a normalisation procedure to account for differences in the recording condition
867 of the different registers (due to divergences in the mouth-to-microphone distance between recordings for
868 instance). Another comment should be made on the importance of using an appropriate feature space in
869 concordance to the size of the dataset. There are some papers reporting a large number of features but using

870 a small dataset. This is certainly discouraged.

871 A final comment should be made in regards to the clinical assessment of the systems presented in lit-
872 erature. The systematic review has demonstrated a lack of validation of the proposed systems in clinical
873 settings where these automatic tools have served to guide or improve the diagnosis of voice impairments.
874 Indeed, most of these systems have been tested under very restricted settings, depending on the dataset that
875 has been used for training, the recording conditions, etc.; factors which might hinder the generalisability of
876 the results. It is necessary to test out the validity of the AVCA methodology in more realistic scenarios,
877 where its ability to contribute to the diagnosis of voice pathologies is tested.

878 To finalise this paper, Table 1 presents, in the authors' opinion, a list of some interesting works presenting
879 AVCA systems, in the hope that they might result useful for new comers to the field. In the second part of
880 this series, entitled "On the design of automatic voice condition analysis systems. Part II: review of speaker
881 recognition techniques and study on the effects of different variability factors" we will introduce a series of
882 experiments following the methodologies described in this paper, using diverse corpora and analysing the
883 effects of certain variability factors in the design of AVCA systems.

Table 1. Some relevant works in literature according to the taxonomy presented in Figure 3. RS.: running speech; AC.: acoustic and temporal features; SC.: spectral/cepstral features; Pn.: perturbation features; Cx.: complexity features; Ot.: other type of features.

Authors	Material	Database	Characterization					Dimensionality reduction	Decision making	Validation	Results
			AC.	SC.	Pn.	Cx.	Ot.				
[156]	/a/	MEEI						–	GMM	Crossvalidation	ACC=98%
[75]	RS.	privative						–	Regressor	Boothstrapping	AUC=0.85
[96]	/a/	privative						–	SVM; KNN; committee	Leave-one-out	ACC=95%(detection); ACC=85%(identification)
[146]	/a/+RS.	MEEI						–	ANN	Crossvalidation	ACC=94%(vowel); ACC=96%(RS.)
[155]	/a/	3 databases						–	XGBoost; ANN; isolation forest	Crossvalidation	ACC = 62-73%
[224]	/a/	MEEI						–	QDA	Bootstrapping	ACC=92%
[119]	/a/	MEEI						Max-relevance	SVM	Crossvalidation	ACC=94%(detection); ACC=85-92%(identification)
[94]	/a/	MEEI; SVD; privative						Mann-Whitney U test	SVM; random forest		ACC=68-100%
[23]	/a/+RS.	MEEI						–	LDA	Leave-one-out	ACC=96%(vowel); ACC=96%(RS.)
[43]	/a/	MEEI						–	ANN	Crossvalidation	ACC=90%
[213]	/a/	MEEI						–	LDA	Leave-one-out	ACC=96%
[63]	/a/	privative						wrapper-based feature selection	SVM; random forests	Cross-validation	ACC=87%
[90]	/a/	privative						GA	SVM	Random-split	ACC=98%(multimodal)
[191, 192]	/a/	privative						4 filter-based	SVM; random forests	Cross-validation	ACC=98%

884 **Acknowledgment**

885 This work was supported by the Ministry of Economy and Competitiveness of Spain under grant DPI2017-
886 83405-R1.

887 **References**

- 888 [1] J. Laver, Principles of Phonetics, Cambridge University Press, Cambridge, 1994.
- 889 [2] H. Traunmüller, Conventional, Biological and Environmental Factors in Speech Communication: A
890 Modulation Theory, *Phonetica* 51 (1-3) (1994) 170–183.
- 891 [3] H. Traunmüller, Evidence for demodulation in speech perception, 2000, pp. 790–793.
- 892 [4] T. Kinnunen, H. Li, An overview of text-independent speaker recognition: From features to supervectors,
893 *Speech Communication* 52 (1) (2010) 12–40.
- 894 [5] M. H. Bahari, M. McLaren, H. Van hamme, D. a. Van Leeuwen, Speaker age estimation using i-vectors,
895 *Engineering Applications of Artificial Intelligence* 34 (2014) 99–108.
- 896 [6] D. Childers, K. Wu, Gender recognition from speech. Part II: Fine analysis., *The Journal of the*
897 *Acoustical Society of America* 90 (4 Pt 1) (1991) 1841–56.
- 898 [7] M. El Ayadi, M. S. Kamel, F. Karray, Survey on speech emotion recognition: Features, classification
899 schemes, and databases, *Pattern Recognition* 44 (3) (2011) 572–587.
- 900 [8] B. Schuller, G. Rigoll, Recognising Interest in Conversational Speech – Comparing Bag of Frames
901 and Supra-segmental Features, in: *INTERSPEECH*, 2009, pp. 1999–2002.
- 902 [9] F. Biadsy, *Automatic Dialect and Accent Recognition and its Application to Speech Recognition* (2011).
- 903 [10] J. Benesty, M. M. Sondhi, Y. Huang, *Springer Handbook of Speech Processing*, Springer Berlin Hei-
904 delberg, Berlin, Heidelberg, 2008.
- 905 [11] C. R. Watts, S. N. Awan, Use of Spectral/Cepstral Analyses for Differentiating Normal From Hypo-
906 functional Voices in Sustained Vowel and Continuous Speech Contexts, *Journal of Speech, Language,*
907 *and Hearing Research* 54 (December) (2011) 1525–1538.
- 908 [12] M. Anniko, M. Bernal-Sprekelsen, V. Bonkowsky, P. Bradley, S. Iurato, *Otorhinolaryngology, Head*
909 *and Neck Surgery*, Springer, 2010.
- 910 [13] S. N. Awan, N. Roy, S. M. Cohen, Exploring the relationship between spectral and cepstral measures
911 of voice and the voice handicap index (VHI), *Journal of Voice* 28 (4) (2014) 430–439.

- 912 [14] J. K. Casper, R. Leonard, Understanding Voice Problems: A Physiological Perspective for Diagnosis
913 and Treatment, illustrate Edition, Lippincott Williams & Wilkins, 2006.
- 914 [15] S. N. Awan, N. Roy, C. Dromey, Estimating dysphonia severity in continuous speech: application of a
915 multi-parameter spectral/cepstral model., *Clinical linguistics & phonetics* 23 (11) (2009) 825–841.
- 916 [16] J. I. Godino-Llorente, N. Sáenz-Lechón, V. Oasma-Ruiz, S. Aguilera-Navarro, P. Gómez-Vilda, An
917 integrated tool for the diagnosis of voice disorders, *Medical Engineering & Physics* 28 (3) (2006) 276–
918 289.
- 919 [17] A. S.-L.-H. Association, Definitions of Communication Disorders and Variations, Tech. rep. (1993).
- 920 [18] J. B. Snow, J. J. Ballenger, Ballenger’s Otorhinolaryngology Head and Neck Surgery, 2003.
- 921 [19] N. B. Anderson, G. H. Shames, Human Communication Disorders: An Introduction, 8th Edition,
922 Pearson, 2010.
- 923 [20] A. E. Aronson, D. Bless, Clinical Voice Disorders, 4th Edition, Thieme, 2009.
- 924 [21] A. Sprecher, A. Olszewski, J. J. Jiang, Y. Zhang, Updating signal typing in voice: addition of type 4
925 signals., *The Journal of the Acoustical Society of America* 127 (6) (2010) 3710–16.
- 926 [22] I. Titze, Workshop on Acoustic Voice Analysis, National Centre for Voice and Speech, America (1994)
927 1–36.
- 928 [23] V. Parsa, D. G. Jamieson, Acoustic discrimination of pathological voice: sustained vowels versus
929 continuous speech, *Journal of speech, language, and hearing research : JSLHR* 44 (2) (2001) 327.
- 930 [24] J. Schoentgen, Spectral models of additive and modulation noise in speech and phonatory excitation
931 signals, *The Journal of the Acoustical Society of America* 113 (1) (2003) 553.
- 932 [25] A. Alpan, Y. Maryn, A. Kacha, F. Grenez, J. Schoentgen, Multi-band dysperiodicity analyses of
933 disordered connected speech, *Speech Communication* 53 (1) (2011) 131–141.
- 934 [26] J. Hanquinet, F. Grenez, J. Schoentgen, Synthesis of Disordered Voices, in: International Conference
935 on Non-Linear Speech Processing, NOLISP 2005, 2006, pp. 231–241.
- 936 [27] M. Vasilakis, Y. Stylianou, Spectral jitter modeling and estimation, *Biomedical Signal Processing and*
937 *Control* 4 (3) (2009) 183–193.
- 938 [28] T. L. Eadie, P. C. Doyle, Classification of Dysphonic Voice: Acoustic and Auditory-Perceptual Mea-
939 sures, *Journal of Voice* 19 (1) (2005) 1–14.
- 940 [29] L. Baghai-Ravary, S. W. Beet, Automatic Speech Signal Analysis for Clinical Diagnosis and Assessment
941 of Speech Disorders, *Automatic Speech Signal Analysis for Clinical Diagnosis and Assessment of Speech*
942 *Disorders (Fuchs 2005)* (2013) 7–12.

- 943 [30] J. E. Huber, E. T. Stathopoulos, G. M. Curione, T. A. Ash, K. Johnson, Formants of children, women,
944 and men: the effects of vocal intensity variation., *The Journal of the Acoustical Society of America*
945 106 (3 Pt 1) (1999) 1532–42.
- 946 [31] C. Moers, B. Möbius, F. Rosanowski, E. Nöth, U. Eysholdt, T. Haderlein, Vowel- and text-based
947 cepstral analysis of chronic hoarseness, *Journal of Voice* 26 (4) (2012) 416–424.
- 948 [32] M. Putzer, W. J. Barry, Instrumental dimensioning of normal and pathological phonation using acoustic
949 measurements., *Clinical linguistics & phonetics* 22 (6) (2008) 407–20.
- 950 [33] R. D. Kent, *The MIT Encyclopedia of Communication Disorders*, MIT Press, 2004.
- 951 [34] B. H. Jacobson, A. Johnson, C. Grywalski, A. Silbergleit, G. Jacobson, M. S. Benninger, C. W.
952 Newman, The voice handicap index (VHI), *American Journal of Speech-Language Pathology* 6 (3)
953 (1997) 66.
- 954 [35] G. B. Kempster, B. R. Gerratt, K. Verdolini Abbott, J. Barkmeier-Kraemer, R. E. Hillman, Consensus
955 Auditory-Perceptual Evaluation of Voice: Development of a Standardized Clinical Protocol, *American*
956 *Journal of Speech-Language Pathology* 18 (2) (2009) 124.
- 957 [36] A. S.-L.-H. Association, [Online; accessed 18-August-2016].
- 958 [37] R. Fraile, J. I. Godino-Llorente, N. Sáenz-Lechón, V. Osma-Ruiz, J. M. Gutiérrez-Arriola, Character-
959 ization of Dysphonic Voices by Means of a Filterbank-Based Spectral Analysis: Sustained Vowels and
960 Running Speech, *Journal of Voice* 27 (1) (2013) 11–23.
- 961 [38] I. R. Titze, W. S. Winholtz, Effect of microphone type and placement on voice perturbation measure-
962 ments, *Journal of Speech & Hearing Research* 36 (6) (1993) 1177–1190.
- 963 [39] J. G. Svec, S. Granqvist, Guidelines for Selecting Microphones for Human Voice Production Research,
964 *American Journal of Speech-Language Pathology* 19 (4) (2010) 356–368.
- 965 [40] A. Ghio, G. Pouchoulin, B. Teston, S. Pinto, C. Fredouille, C. De Looze, D. Robert, F. Viallet,
966 A. Giovanni, How to manage sound, physiological and clinical data of 2500 dysphonic and dysarthric
967 speakers?, *Speech Communication* 54 (5) (2012) 664–679.
- 968 [41] Massachusetts Eye and Ear Infirmary, *Voice disorders database, version.1.03 [cd-rom]*, Lincoln Park,
969 NJ: Kay Elemetrics Corp (1994).
- 970 [42] V. Parsa, D. G. Jamieson, Identification of pathological voices using glottal noise measures, *Journal of*
971 *speech, language, and hearing research* 43 (2) (2000) 469.
- 972 [43] N. Sáenz-Lechón, J. I. Godino-Llorente, V. Osma-Ruiz, P. Gómez-Vilda, Methodological issues in the
973 development of automatic systems for voice pathology detection, *Biomedical Signal Processing and*
974 *Control* 1 (2) (2006) 120–128.

- 975 [44] N. Malyska, T. Quatieri, D. Sturim, Automatic Dysphonia Recognition using Biologically-Inspired
976 Amplitude-Modulation Features, in: Proceedings. (ICASSP '05). IEEE International Conference on
977 Acoustics, Speech, and Signal Processing, 2005., Vol. 1, IEEE, 2005, pp. 873–876.
- 978 [45] Saarbrücken voice database.
979 URL <http://www.stimmdatenbank.coli.uni-saarland.de/index.php4>
- 980 [46] J. I. Godino-Llorente, V. Osma-Ruiz, N. Sáenz-Lechón, I. Cobeta-Marco, R. González-Herranz,
981 C. Ramírez-Calvo, Acoustic analysis of voice using WPCVox: a comparative study with Multi Di-
982 mensional Voice Program, *European Archives of Oto-Rhino-Laryngology* 265 (4) (2008) 465–476.
- 983 [47] A. Al-Nasheri, G. Muhammad, M. Alsulaiman, Z. Ali, T. a. Mesallam, M. Farahat, K. H. Malki,
984 M. A. Bencherif, An Investigation of Multidimensional Voice Program Parameters in Three Different
985 Databases for Voice Pathology Detection and Classification, *Journal of Voice* 31 (1) (2017) 113.e9–
986 113.e18.
- 987 [48] T. A. Mesallam, M. Farahat, K. H. Malki, M. Alsulaiman, Z. Ali, A. Al-nasheri, G. Muhammad,
988 Development of the Arabic Voice Pathology Database and Its Evaluation by Using Speech Features
989 and Machine Learning Algorithms, *Journal of Healthcare Engineering* 2017 (2017) 1–13.
- 990 [49] C. D. P. Crovato, A. Schuck, The Use of Wavelet Packet Transform and Artificial Neural Networks in
991 Analysis and Classification of Dysphonic Voices, *IEEE Transactions on Biomedical Engineering* 54 (10)
992 (2007) 1898–1900.
- 993 [50] R. Behroozmand, F. Almasganj, Optimal selection of wavelet-packet-based features using genetic al-
994 gorithm in pathological assessment of patients' speech signal with unilateral vocal fold paralysis, *Com-
995 puters in Biology and Medicine* 37 (4) (2007) 474–485.
- 996 [51] C. Fabris, W. De Colle, G. Sparacino, Voice disorders assessed by (cross-) Sample Entropy of elec-
997 troglottogram and microphone signals, *Biomedical Signal Processing and Control* 8 (6) (2013) 920–926.
- 998 [52] L. A. Forero M., M. Kohler, M. M. Vellasco, E. Cataldo, Analysis and Classification of Voice Pathologies
999 Using Glottal Signal Parameters, *Journal of Voice* 30 (5) (2016) 549–556.
- 1000 [53] M. Fröhlich, D. Michaelis, H. W. Strube, E. Kruse, Acoustic voice analysis by means of the hoarseness
1001 diagram., *Journal of speech, language, and hearing research : JSLHR* 43 (3) (2000) 706–20.
- 1002 [54] A. Giovanni, M. Ouaknine, J. M. Triglia, Determination of largest Lyapunov exponents of vocal signal:
1003 application to unilateral laryngeal paralysis., *Journal of voice : official journal of the Voice Foundation*
1004 13 (3) (1999) 341–54.
- 1005 [55] S. Hadjitodorov, B. Boyanov, B. Teston, Laryngeal pathology detection by means of class-specific
1006 neural maps, *IEEE Transactions on Information Technology in Biomedicine* 4 (1) (2000) 68–73.

- 1007 [56] J. Jiang, Y. Zhang, Nonlinear dynamic analysis of speech from pathological subjects, *Electronics*
1008 *Letters* 38 (6) (2002) 294–295.
- 1009 [57] A novel hybrid of genetic algorithm and ANN for developing a high efficient method for vocal fold
1010 pathology diagnosis, *EURASIP Journal on Audio, Speech, and Music Processing* 2015 (1) (2015) 3.
- 1011 [58] P. Mitev, S. Hadjitodorov, A method for turbulent noise estimation in voiced signals, *Med. Biol. Eng.*
1012 *Comput* 38 (2000) 625–631.
- 1013 [59] M. Petrović-Lazić, S. Babac, M. Vuković, R. Kosanović, Z. Ivanković, Acoustic Voice Analysis of
1014 Patients With Vocal Fold Polyp, *Journal of Voice* 25 (1) (2011) 94–97.
- 1015 [60] P. Saidi, F. Almasganj, Voice Disorder Signal Classification Using M-Band Wavelets and Support
1016 Vector Machine, *Circuits, Systems, and Signal Processing* 34 (2015) 2727–2738.
- 1017 [61] P. R. Scalassara, M. E. Dajer, C. D. Maciel, R. C. Guido, J. C. Pereira, Relative entropy measures
1018 applied to healthy and pathological voice characterization, *Applied Mathematics and Computation*
1019 207 (1) (2009) 95–108.
- 1020 [62] V. Uloza, A. Verikas, M. Bacauskiene, A. Gelzinis, R. Pribuisiene, M. Kasetas, V. Saferis, Categorizing
1021 normal and pathological voices: Automated and perceptual categorization, *Journal of Voice* 25 (6)
1022 (2011) 700–708.
- 1023 [63] E. Vaiciukynas, A. Verikas, A. Gelzinis, M. Bacauskiene, Z. Kons, A. Satt, R. Hoory, Fusion of voice
1024 signal information for detection of mild laryngeal pathology, *Applied Soft Computing* 18 (2014) 91–103.
- 1025 [64] J. H. Van Stan, D. D. Mehta, S. M. Zeitels, J. A. Burns, A. M. Barbu, R. E. Hillman, Average
1026 ambulatory measures of sound pressure level, fundamental frequency, and vocal dose do not differ
1027 between adult females with phonotraumatic lesions and matched control subjects, *Annals of Otolaryngology,*
1028 *Rhinology and Laryngology* 124 (11) (2015) 864–874.
- 1029 [65] A. Verikas, A. Gelzinis, E. Vaiciukynas, M. Bacauskiene, J. Minelga, M. Hållander, V. Uloza,
1030 E. Padervinskis, Data dependent random forest applied to screening for laryngeal disorders through
1031 analysis of sustained phonation: Acoustic versus contact microphone, *Medical Engineering & Physics*
1032 37 (2) (2015) 210–218.
- 1033 [66] M. N. Vieira, F. R. McInnes, M. a. Jack, On the influence of laryngeal pathologies on acoustic and
1034 electroglottographic jitter measures., *The Journal of the Acoustical Society of America* 111 (2) (2002)
1035 1045–1055.
- 1036 [67] Y. Zhang, J. J. Jiang, Nonlinear dynamic analysis in signal typing of pathological human voices,
1037 *Electronics Letters* 39 (13) (2003) 1021.

- 1038 [68] P. Henriquez, J. B. Alonso-Hernandez, M. A. Ferrer, C. M. Travieso, J. I. Godino-Llorente, F. Diaz-de
1039 Maria, Characterization of Healthy and Pathological Voice Through Measures Based on Nonlinear
1040 Dynamics, *IEEE Transactions on Audio, Speech, and Language Processing* 17 (6) (2009) 1186–1195.
- 1041 [69] M. de Oliveira Rosa, J. Pereira, M. Grellet, Adaptive estimation of residue signal for voice pathology
1042 diagnosis, *IEEE Transactions on Biomedical Engineering* 47 (1) (2000) 96–104.
- 1043 [70] L. Landini, C. Manfredi, V. Positano, M. Santarelli, N. Vanello, Non-linear prediction for oesophageal
1044 voice analysis, *Medical Engineering & Physics* 24 (7-8) (2002) 529–533.
- 1045 [71] C. M. Travieso, J. B. Alonso, J. Orozco-Arroyave, J. Vargas-Bonilla, E. Nöth, A. G. Ravelo-García,
1046 Detection of different voice diseases based on the nonlinear characterization of speech signals, *Expert
1047 Systems with Applications* 82 (2017) 184–195.
- 1048 [72] V. Péan, M. Ouayoun, B. C. Fugain, A Fractal Approach to Normal and Pathological Voices, *Acta
1049 Oto-Laryngologica* 120 (2) (2000) 222–224.
- 1050 [73] A. Al-nasheri, G. Muhammad, M. Alsulaiman, Z. Ali, Investigation of Voice Pathology Detection and
1051 Classification on Different Frequency Regions Using Correlation Functions, *Journal of Voice* 31 (1)
1052 (2017) 3–15.
- 1053 [74] V. Majidnezhad, A HTK-based Method for Detecting Vocal Fold Pathology, *Acta Informatica Medica
1054* 22 (4) (2014) 246.
- 1055 [75] S. N. Awan, N. Roy, D. Zhang, S. M. Cohen, Validation of the Cepstral Spectral Index of Dysphonia
1056 (CSID) as a Screening Tool for Voice Disorders: Development of Clinical Cutoff Scores, *Journal of
1057 Voice* 30 (2) (2016) 130–144.
- 1058 [76] S. Y. Lowell, R. H. Colton, R. T. Kelley, Y. C. Hahn, Spectral- and Cepstral-Based Measures During
1059 Continuous Speech: Capacity to Distinguish Dysphonia and Consistency Within a Speaker, *Journal
1060 of Voice* 25 (5) (2010) e223–e232.
- 1061 [77] D. D. Mehta, J. H. Van Stan, M. Zanartu, M. Ghassemi, J. V. Guttag, V. M. Espinoza, J. P. Cortes,
1062 H. A. n. Cheyne, R. E. Hillman, Using Ambulatory Voice Monitoring to Investigate Common Voice
1063 Disorders: Research Update., *Frontiers in bioengineering and biotechnology* 3 (October) (2015) 155.
- 1064 [78] E. Ma, J. Robertson, C. Radford, S. Vagne, R. El-Halabi, E. Yiu, Reliability of Speaking and Maximum
1065 Voice Range Measures in Screening for Dysphonia, *Journal of Voice* 21 (4) (2007) 397–406.
- 1066 [79] M. A. Little, D. A. Costello, M. L. Harries, Objective Dysphonia Quantification in Vocal Fold Paralysis:
1067 Comparing Nonlinear With Classical Measures, *Journal of Voice* 25 (1) (2011) 21–31.

- 1068 [80] L. F. Brinca, A. P. F. Batista, A. I. Tavares, I. C. Gonçalves, M. L. Moreno, Use of Cepstral Analyses
1069 for Differentiating Normal From Dysphonic Voices: A Comparative Study of Connected Speech Versus
1070 Sustained Vowel in European Portuguese Female Speakers, *Journal of Voice* 28 (3) (2014) 282–286.
- 1071 [81] C. Fredouille, G. Pouchoulin, J.-F. Bonastre, M. Azzarello, A. Giovanni, A. Ghio, Application of
1072 Automatic Speaker Recognition techniques to pathological voice assessment (dysphonia), in: *Proceed-*
1073 *ings of European Conference on Speech Communication and Technology (Eurospeech 2005)*, 2005, pp.
1074 149–152.
- 1075 [82] F. Klingholtz, Acoustic recognition of voice disorders: a comparative study of running speech versus
1076 sustained vowels., *The Journal of the Acoustical Society of America* 87 (5) (1990) 2218–2224.
- 1077 [83] J.-W. Lee, S. Kim, H.-G. Kang, Detecting pathological speech using contour modeling of harmonic-
1078 to-noise ratio, in: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing*
1079 *(ICASSP)*, IEEE, 2014, pp. 5969–5973.
- 1080 [84] T. Lee, Y. Liu, P.-W. Huang, J.-T. Chien, W. K. Lam, Y. T. Yeung, T. K. T. Law, K. Y. S. Lee,
1081 A. P.-H. Kong, S.-P. Law, Automatic speech recognition for acoustical analysis and assessment of
1082 cantonese pathological voice and speech, in: *2016 IEEE International Conference on Acoustics, Speech*
1083 *and Signal Processing (ICASSP)*, Vol. 2016-May, IEEE, 2016, pp. 6475–6479.
- 1084 [85] Y. Maryn, P. Corthals, P. Van Cauwenberge, N. Roy, M. De Bodt, Toward Improved Ecological Validity
1085 in the Acoustic Measurement of Overall Voice Quality: Combining Continuous Speech and Sustained
1086 Vowels, *Journal of Voice* 24 (5) (2010) 540–555.
- 1087 [86] C. Middag, Y. Saeys, J.-p. Martens, Towards an ASR-free objective analysis of pathological speech,
1088 in: *Interspeech*, no. September, 2010, pp. 294–297.
- 1089 [87] K. Manickam, C. Moore, T. Willard, N. Slevin, Quantifying aberrant phonation using approximate
1090 entropy in electrolaryngography, *Speech Communication* 47 (3) (2005) 312–321.
- 1091 [88] C. Moore, K. Manickam, T. Willard, S. Jones, N. Slevin, S. Shalet, Spectral pattern complexity
1092 analysis and the quantification of voice normality in healthy and radiotherapy patient groups, *Medical*
1093 *Engineering & Physics* 26 (4) (2004) 291–301.
- 1094 [89] E. Vaiciukynas, A. Verikas, A. Gelzinis, M. Bacauskiene, J. Minelga, M. Hållander, E. Padervinskis,
1095 V. Uloza, Fusing voice and query data for non-invasive detection of laryngeal disorders, *Expert Systems*
1096 *with Applications* 42 (22) (2015) 8445–8453.
- 1097 [90] A. Verikas, A. Gelzinis, M. Bacauskiene, M. Hållander, V. Uloza, M. Kaseta, Combining image, voice,
1098 and the patient’s questionnaire data to categorize laryngeal disorders, *Artificial Intelligence in Medicine*
1099 49 (1) (2010) 43–50.

- 1100 [91] J. W. Lee, H.-G. Kang, J.-Y. Choi, Y.-I. Son, An Investigation of Vocal Tract Characteristics for
1101 Acoustic Discrimination of Pathological Voices, *BioMed Research International* 2013 (2013) 1–11.
- 1102 [92] J. R. Orozco, J. D. Arias, J. F. Vargas, M. C. González Rátiva, E. Nöth, New Spanish speech corpus
1103 database for the analysis of people suffering from Parkinson’s disease, *LREC 2014. Proceedings of the*
1104 *Ninth International Conference on Language Resources and Evaluation* (2014) 342–347.
- 1105 [93] A. Maier, T. Haderlein, U. Eysholdt, F. Rosanowski, A. Batliner, M. Schuster, E. Nöth, PEAKS -
1106 A system for the automatic evaluation of voice and speech disorders, *Speech Communication* 51 (5)
1107 (2009) 425–437.
- 1108 [94] J. Mekyska, E. Janousova, P. Gómez-Vilda, Z. Smekal, I. Rektorova, I. Eliasova, M. Kostalova,
1109 M. Mrackova, J. B. Alonso-Hernandez, M. Faundez-Zanuy, K. López-de Ipiña, Robust and complex
1110 approach of pathological speech signal analysis, *Neurocomputing* 167 (2015) 94–111.
- 1111 [95] V. Uloza, V. Saferis, I. Uloziene, Perceptual and Acoustic Assessment of Voice Pathology and the
1112 Efficacy of Endolaryngeal Phonomicrosurgery, *Journal of Voice* 19 (1) (2005) 138–145.
- 1113 [96] A. Gelzinis, A. Verikas, M. Bacauskiene, Automated speech analysis applied to laryngeal disease cat-
1114 egorization, *Computer Methods and Programs in Biomedicine* 91 (1) (2008) 36–47.
- 1115 [97] B. Radish Kumar, J. S. Bhat, N. Prasad, Cepstral analysis of voice in persons with vocal nodules,
1116 *Journal of Voice* 24 (6) (2010) 651–653.
- 1117 [98] Y. Zhang, C. McGilligan, L. Zhou, M. Vig, J. J. Jiang, Nonlinear dynamic analysis of voices before
1118 and after surgical excision of vocal polyps., *The Journal of the Acoustical Society of America* 115 (5
1119 Pt 1) (2004) 2270–7.
- 1120 [99] A. Roviroso, E. Martínez-Celdrán, A. Ortega, C. Ascaso, R. Abellana, M. Velasco, M. Bonet, C. Herrera,
1121 F. Casas, R. M. Francisco, M. Arenas, V. Hernández, A. Sánchez-Reyes, C. León, J. Traserra, A. Biete,
1122 Acoustic analysis after radiotherapy in T1 vocal cord carcinoma: a new approach to the analysis of
1123 voice quality, *International Journal of Radiation Oncology*Biology*Physics* 47 (1) (2000) 73–79.
- 1124 [100] T. Ritchings, M. McGillion, C. Moore, Pathological voice quality assessment using artificial neural
1125 networks, *Medical Engineering & Physics* 24 (7-8) (2002) 561–564.
- 1126 [101] P. Aichinger, I. Roesner, M. Leonhard, B. Schneider-Stickler, D.-M. Denk-Linnert, W. Bigenzahn,
1127 A. K. Fuchs, M. Hagmüller, G. Kubin, Comparison of an audio-based and a video-based approach for
1128 detecting diplophonia, *Biomedical Signal Processing and Control* 31 (31) (2017) 576–585.
- 1129 [102] G. Schlotthauer, M. E. Torres, M. C. Jackson-Menaldi, A Pattern Recognition Approach to Spasmodic
1130 Dysphonia and Muscle Tension Dysphonia Automatic Classification, *Journal of Voice* 24 (3) (2010)
1131 346–353.

- 1132 [103] Y. Zhang, J. J. Jiang, L. Biazzo, M. Jorgensen, Perturbation and nonlinear dynamic analyses of voices
1133 from patients with unilateral laryngeal paralysis, *Journal of Voice* 19 (4) (2005) 519–528.
- 1134 [104] E. Goldshtein, A. Tarasiuk, Y. Zigel, Automatic detection of obstructive sleep apnea using speech
1135 signals., *IEEE transactions on bio-medical engineering* 58 (5) (2011) 1373–82.
- 1136 [105] J. L. Blanco-Murillo, L. A. Hernández, R. Fernández-Pozo, D. Ramos, Improving Automatic Detection
1137 of Obstructive Sleep Apnea Through Nonlinear Analysis of Sustained Speech, *Cognitive Computation*.
- 1138 [106] P. Vijayalakshmi, M. R. Reddy, D. O’Shaughnessy, Acoustic Analysis and Detection of Hypernasality
1139 Using a Group Delay Function, *IEEE Transactions on Biomedical Engineering* 54 (4) (2007) 621–629.
- 1140 [107] M. Little, P. McSharry, E. Hunter, J. Spielman, L. Ramig, Suitability of Dysphonia Measurements for
1141 Telemonitoring of Parkinson’s Disease, *IEEE Transactions on Biomedical Engineering* 56 (4) (2009)
1142 1015–1022.
- 1143 [108] J. Rusz, R. Cmejla, T. Tykalova, H. Ruzickova, J. Klempir, V. Majerova, J. Picmausova, J. Roth,
1144 E. Ruzicka, Imprecise vowel articulation as a potential early marker of Parkinson’s disease: effect of
1145 speaking task., *The Journal of the Acoustical Society of America* 134 (3) (2013) 2171–81.
- 1146 [109] K. López-de Ipiña, P. Calvo, M. Faundez-Zanuy, P. Clavé, W. Nascimento, U. Martinez de Lizarduy,
1147 D. Alvarez, V. Arreola, O. Ortega, J. Mekyska, P. Sanz-Cartagena, Automatic voice analysis for
1148 dysphagia detection, *Speech, Language and Hearing* 21 (2) (2018) 86–89.
- 1149 [110] J. S. Ryu, S. R. Park, K. H. Choi, Prediction of laryngeal aspiration using voice analysis, *American*
1150 *Journal of Physical Medicine and Rehabilitation* 83 (10) (2004) 753–757.
- 1151 [111] M. S. F. C. de Macedo, K. M. Costa, M. da Silva Filho, Voice disorder in systemic lupus erythematosus,
1152 *PLOS ONE* 12 (4) (2017) e0175893.
- 1153 [112] K. Paliwal, K. Wojcicki, Effect of analysis window duration on speech intelligibility, *IEEE Signal*
1154 *Processing Letters* 15 (2008) 785–788.
- 1155 [113] X. Huang, A. Acero, H. W. Hon, *Spoken Language Processing: A Guide to Theory, Algorithm, and*
1156 *System Development*, Prentice Hall PTR, 2001.
- 1157 [114] Ö. Eskidere, A. GürhanlÄs, Voice Disorder Classification Based on Multitaper Mel Frequency Cepstral
1158 Coefficients Features, *Computational and Mathematical Methods in Medicine* 2015 (2015) 1–12.
- 1159 [115] A. A. Dibazar, T. W. Berger, S. S. Narayanan, Pathological Voice Assessment, in: *2006 International*
1160 *Conference of the IEEE Engineering in Medicine and Biology Society*, IEEE, 2006, pp. 1669–1673.
- 1161 [116] A. Alpan, J. Schoentgen, Y. Maryn, F. Grenez, Automatic perceptual categorization of disordered
1162 connected speech, in: *INTERSPEECH 2010, 11th Annual Conference of the International Speech*
1163 *Communication Association*, 2010, pp. 2574–2577.

- 1164 [117] K. Shama, A. Krishna, N. U. Cholayya, Study of Harmonics-to-Noise Ratio and Critical-Band Energy
1165 Spectrum of Speech as Acoustic Indicators of Laryngeal and Voice Pathology, *EURASIP Journal on*
1166 *Advances in Signal Processing* 2007 (1) (2007) 085286.
- 1167 [118] C. Manfredi, M. D’Aniello, P. Bruscaioni, A. Ismaelli, A comparative analysis of fundamental fre-
1168 quency estimation methods with application to pathological voices, *Medical Engineering & Physics*
1169 22 (2) (2000) 135–147.
- 1170 [119] M. Markaki, Y. Stylianou, Voice Pathology Detection and Discrimination Based on Modulation Spec-
1171 tral Features, *IEEE Transactions on Audio, Speech, and Language Processing* 19 (7) (2011) 1938–1948.
- 1172 [120] G. Muhammad, M. Alsulaiman, Z. Ali, T. A. Mesallam, M. Farahat, K. H. Malki, A. Al-nasheri, M. A.
1173 Bencherif, Voice pathology detection using interlaced derivative pattern on glottal source excitation,
1174 *Biomedical Signal Processing and Control* 31 (2017) 156–164.
- 1175 [121] S. Jothilakshmi, Automatic system to detect the type of voice pathology, *Applied Soft Computing* 21
1176 (2014) 244–249.
- 1177 [122] S. Hadjitodorov, P. Mitev, A computer system for acoustic analysis of pathological voices and laryngeal
1178 diseases screening, *Medical Engineering & Physics* 24 (6) (2002) 419–429.
- 1179 [123] M. K. Arjmandi, M. Pooyan, An optimum algorithm in pathological voice quality assessment using
1180 wavelet-packet-based features, linear discriminant analysis and support vector machine, *Biomedical*
1181 *Signal Processing and Control* 7 (1) (2012) 3–19.
- 1182 [124] D. Panek, A. Skalski, J. Gajda, Quantification of Linear and Non-linear Acoustic Analysis Applied to
1183 Voice Pathology Detection, in: E. Piętko, J. Kawa, W. Wieclawek (Eds.), *Information Technologies*
1184 *in Biomedicine*, Vol. 284 of *Advances in Intelligent Systems and Computing*, Springer International
1185 Publishing, Cham, 2014.
- 1186 [125] Y. Qi, R. E. Hillman, C. Milstein, The estimation of signal-to-noise ratio in continuous speech for
1187 disordered voices, *Journal of the Acoustical Society of America* 105 (4) (1999) 2532–2535.
- 1188 [126] Z. Ali, I. Elamvazuthi, M. Alsulaiman, G. Muhammad, Automatic Voice Pathology Detection With
1189 Running Speech by Using Estimation of Auditory Spectrum and Cepstral Coefficients Based on the
1190 All-Pole Model, *Journal of Voice* 30 (6) (2016) 757.e7–757.e19.
- 1191 [127] D. Panek, A. Skalski, J. Gajda, R. Tadeusiewicz, Acoustic Analysis Assessment in Speech Pathology
1192 Detection, *International Journal of Applied Mathematics and Computer Science* 25 (2015) 631–643.
- 1193 [128] M. S. Hossain, Cloud-Supported Cyber-Physical Localization Framework for Patients Monitoring,
1194 *IEEE Systems Journal* 11 (1) (2017) 118–127.

- 1195 [129] J. I. Godino-Llorente, P. Gómez-Vilda, M. Blanco-Velasco, Dimensionality reduction of a pathological
1196 voice quality assessment system based on Gaussian mixture models and short-term cepstral paramete-
1197 ters., *IEEE transactions on bio-medical engineering* 53 (10) (2006) 1943–1953.
- 1198 [130] E. Vaiciukynas, A. Verikas, A. Gelzinis, M. Bacauskiene, V. Uloza, Exploring similarity-based classifi-
1199 cation of larynx disorders from human voice, *Speech Communication* 54 (5) (2012) 601–610.
- 1200 [131] R. Moran, R. Reilly, P. de Chazal, P. Lacy, Telephony-Based Voice Pathology Assessment Using
1201 Automated Speech Analysis, *IEEE Transactions on Biomedical Engineering* 53 (3) (2006) 468–477.
- 1202 [132] G. Muhammad, T. a. Mesallam, K. H. Malki, M. Farahat, A. Mahmood, M. Alsulaiman, Multidirec-
1203 tional regression (MDR)-based features for automatic voice disorder detection, *Journal of Voice* 26 (6)
1204 (2012) 817.e19–817.e27.
- 1205 [133] X. Wang, J. Zhang, Y. Yan, Discrimination Between Pathological and Normal Voices Using GMM-SVM
1206 Approach, *Journal of Voice* 25 (1) (2011) 38–43.
- 1207 [134] M. Hariharan, K. Polat, S. Yaacob, A new feature constituting approach to detection of vocal fold
1208 pathology, *International Journal of Systems Science* 45 (8) (2013) 1622–1634.
- 1209 [135] C. Manfredi, Adaptive noise energy estimation in pathological speech signals, *IEEE Transactions on*
1210 *Biomedical Engineering* 47 (11) (2000) 1538–1543.
- 1211 [136] J. D. Arias-Londoño, J. I. Godino-Llorente, N. Sáenz-Lechón, V. Osma-Ruiz, G. Castellanos-
1212 Domínguez, An improved method for voice pathology detection by means of a HMM-based feature
1213 space transformation, *Pattern Recognition* 43 (9) (2010) 3100–3112.
- 1214 [137] J. B. Alonso-Hernandez, J. De Leon, I. Alonso, M. A. Ferrer, Automatic detection of pathologies in
1215 the voice by HOS based parameters, *Eurasip Journal on Applied Signal Processing* 2001 (4) (2001)
1216 275–284.
- 1217 [138] J. I. Godino-Llorente, V. Osma-Ruiz, N. Sáenz-Lechón, P. Gómez-Vilda, M. Blanco-Velasco, F. Cruz-
1218 Roldán, The Effectiveness of the Glottal to Noise Excitation Ratio for the Screening of Voice Disorders,
1219 *Journal of Voice* 24 (1) (2010) 47–56.
- 1220 [139] A. Dibazar, S. Narayanan, T. Berger, Feature analysis for automatic detection of pathological speech,
1221 in: *Proceedings of the Second Joint 24th Annual Conference and the Annual Fall Meeting of the*
1222 *Biomedical Engineering Society* [Engineering in Medicine and Biology, Vol. 1, IEEE, 2002, pp. 182–
1223 183.
- 1224 [140] J. D. Arias-Londoño, J. I. Godino-Llorente, M. Markaki, Y. Stylianou, On combining information
1225 from modulation spectra and mel-frequency cepstral coefficients for automatic detection of pathological
1226 voices., *Logopedics, phoniatrics, vocology* 36 (2) (2011) 60–69.

- 1227 [141] J. R. Orozco-Arroyave, E. A. Belalcazar-Bolaños, J. D. Arias-Londoño, J. F. Vargas-Bonilla, S. Skodda,
1228 J. Rusz, K. Daqrouq, F. Hönig, E. Nöth, Characterization methods for the detection of multiple voice
1229 disorders: Neurological, functional, and laryngeal diseases, *IEEE Journal of Biomedical and Health*
1230 *Informatics* 19 (6) (2015) 1820–1828.
- 1231 [142] G. Muhammad, M. F. Alhamid, M. Alsulaiman, B. Gupta, Edge Computing with Cloud for Voice
1232 Disorder Assessment and Treatment, *IEEE Communications Magazine* 56 (4) (2018) 60–65.
- 1233 [143] G. Muhammad, S. M. M. Rahman, A. Alelaiwi, A. Alamri, Smart Health Solution Integrating IoT and
1234 Cloud: A Case Study of Voice Pathology Monitoring, *IEEE Communications Magazine* 55 (1) (2017)
1235 69–73.
- 1236 [144] T. Dubuisson, T. Dutoit, B. Gosselin, M. Remacle, On the Use of the Correlation between Acoustic
1237 Descriptors for the Normal/Pathological Voices Discrimination, *EURASIP Journal on Advances in*
1238 *Signal Processing* 2009 (1) (2009) 173967.
- 1239 [145] D. Michaelis, M. Fröhlich, H. W. Strube, Selection and combination of acoustic features for the de-
1240 scription of pathologic voices. 103 (3) (1998) 1628–1639.
- 1241 [146] J. I. Godino-Llorente, R. Fraile, N. Sáenz-Lechón, V. Osma-Ruiz, P. Gómez-Vilda, Automatic detection
1242 of voice impairments from text-dependent running speech, *Biomedical Signal Processing and Control*
1243 4 (3) (2009) 176–182.
- 1244 [147] J. I. Godino-Llorente, P. Gomez-Vilda, Automatic Detection of Voice Impairments by Means of Short-
1245 Term Cepstral Parameters and Neural Network Based Detectors, *IEEE Transactions on Biomedical*
1246 *Engineering* 51 (2) (2004) 380–384.
- 1247 [148] M. Vasilakis, Y. Stylianou, Voice pathology detection based on short-term jitter estimations in running
1248 speech, *Folia Phoniatria et Logopaedica* 61 (3) (2009) 153–170.
- 1249 [149] C. Moore, K. Manickam, N. Slevin, Collective spectral pattern complexity analysis of voicing in normal
1250 males and larynx cancer patients following radiotherapy, *Biomedical Signal Processing and Control*
1251 1 (2) (2006) 113–119.
- 1252 [150] G. Daza-Santacoloma, J. D. Arias-Londono, J. I. Godino-Llorente, N. Saenz-Lechon, V. Osma-Ruiz,
1253 G. Castellanos-Dominguez, Dynamic Feature Extraction: an Application To Voice Pathology Detec-
1254 tion, *Intelligent Automation and Soft Computing* 15 (4) (2009) 667–682.
- 1255 [151] S. H. Fang, Y. Tsao, M. J. Hsiao, J. Y. Chen, Y. H. Lai, F. C. Lin, C. T. Wang, Detection
1256 of Pathological Voice Using Cepstrum Vectors: A Deep Learning Approach, *Journal of Voice-*
1257 *doi:10.1016/j.jvoice.2018.02.003.*

- 1258 [152] R. Fraile, N. Sáenz-Lechón, J. I. Godino-Llorente, V. Osma-Ruiz, C. Fredouille, Automatic Detection of
1259 Laryngeal Pathologies in Records of Sustained Vowels by Means of Mel-Frequency Cepstral Coefficient
1260 Parameters and Differentiation of Patients by Sex, *Folia Phoniatria et Logopaedica* 61 (3) (2009)
1261 146–152.
- 1262 [153] G. Muhammad, G. Altuwaijri, M. Alsulaiman, Z. Ali, T. A. Mesallam, M. Farahat, K. H. Malki, A. Al-
1263 nasheri, Automatic voice pathology detection and classification using vocal tract area irregularity,
1264 *Biocybernetics and Biomedical Engineering* 36 (2) (2016) 309–317.
- 1265 [154] G. Muhammad, M. Melhem, Pathological voice detection and binary classification using MPEG-7
1266 audio features, *Biomedical Signal Processing and Control* 11 (1) (2014) 1–9.
- 1267 [155] P. Harar, Z. Galaz, J. B. Alonso-Hernandez, J. Mekyska, R. Burget, Z. Smekal, Towards robust voice
1268 pathology detection, *Neural Computing and Applications* (2018) 1–11.
- 1269 [156] J. D. Arias-Londoño, J. I. Godino-Llorente, N. Sáenz-Lechón, V. Osma-Ruiz, G. Castellanos-
1270 Domínguez, Automatic detection of pathological voices using complexity measures, noise parameters,
1271 and mel-cepstral coefficients, *IEEE Transactions on Biomedical Engineering* 58 (2) (2011) 370–379.
- 1272 [157] J. I. Godino-Llorente, P. Gómez-Vilda, F. Cruz-Roldán, M. Blanco-Velasco, R. Fraile, Pathological
1273 Likelihood Index as a Measurement of the Degree of Voice Normality and Perceived Hoarseness, *Journal*
1274 *of Voice* 24 (6) (2010) 667–677.
- 1275 [158] D. Martínez González, E. Lleida, A. Ortega, A. Miguel, Score level versus audio level fusion for
1276 voice pathology detection on the Saarbrücken Voice Database, *Communications in Computer and*
1277 *Information Science* 328 CCIS (2012) 110–120.
- 1278 [159] D. Martínez González, E. Lleida, A. Ortega, A. Miguel, J. Villalba, Voice pathology detection on the
1279 Saarbrücken Voice Database with calibration and fusion of scores using multifocal toolkit, *Communi-*
1280 *cations in Computer and Information Science* 328 CCIS (2012) 99–109.
- 1281 [160] G. Muhammad, M. F. Alhamid, M. S. Hossain, A. S. Almogren, A. V. Vasilakos, Enhanced Living by
1282 Assessing Voice Pathology Using a Co-Occurrence Matrix, *Sensors* 17 (2) (2017) 267.
- 1283 [161] M. Hariharan, K. Polat, R. Sindhu, S. Yaacob, A hybrid expert system approach for telemonitoring of
1284 vocal fold pathology, *Applied Soft Computing* 13 (10) (2013) 4148–4161.
- 1285 [162] M. S. Hossain, G. Muhammad, Healthcare Big Data Voice Pathology Assessment Framework, *IEEE*
1286 *Access* 4 (2016) 7806–7815.
- 1287 [163] L. Moro-Velázquez, J. A. Gómez-García, J. I. Godino-Llorente, J. Villalba, J. R. Orozco-Arroyave,
1288 N. Dehak, Analysis of speaker recognition methodologies and the influence of kinetic changes to auto-
1289 matically detect Parkinson’s Disease, *Applied Soft Computing* 62 (2018) 649–666.

- 1290 [164] B. Ghoraani, S. Krishnan, A Joint Time-Frequency and Matrix Decomposition Feature Extraction
1291 Methodology for Pathological Voice Classification, *EURASIP Journal on Advances in Signal Processing*
1292 2009 (2009) 1–12.
- 1293 [165] J. Nayak, P. Bhat, R. Acharya, U. Aithal, Classification and analysis of speech abnormalities, *ITBM-*
1294 *RBM* 26 (5-6) (2005) 319–327.
- 1295 [166] H. Ghasemzadeh, M. Tajik Khass, M. Khalil Arjmandi, M. Pooyan, Detection of vocal disorders based
1296 on phase space parameters and Lyapunov spectrum, *Biomedical Signal Processing and Control* 22
1297 (2015) 135–145.
- 1298 [167] J. J. Jiang, Y. Zhang, C. McGilligan, Chaos in voice, from modeling to measurement, *Journal of Voice*
1299 20 (1) (2006) 2–17.
- 1300 [168] M. Kaleem, B. Ghoraani, A. Guergachi, S. Krishnan, Pathological speech signal analysis and classifi-
1301 cation using empirical mode decomposition., *Medical & biological engineering & computing* 51 (2013)
1302 811–21.
- 1303 [169] H. Cordeiro, J. Fonseca, I. Guimarães, C. Meneses, Hierarchical Classification and System Combina-
1304 tion for Automatically Identifying Physiological and Neuromuscular Laryngeal Pathologies, *Journal of*
1305 *Voice*.
- 1306 [170] T. Drugman, T. Dubuisson, T. Dutoit, Phase-based information for voice pathology detection, in: 2011
1307 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2011,
1308 pp. 4612–4615.
- 1309 [171] T. Drugman, T. Dubuisson, T. Dutoit, On the mutual information between source and filter contri-
1310 butions for voice pathology detection, in: *Proceedings of the Annual Conference of the International*
1311 *Speech Communication Association, INTERSPEECH, 2009*, pp. 1463–1466.
- 1312 [172] P. Gómez-Vilda, R. Fernández-Baillo, A. Nieto, F. Díaz, F. Fernández-Camacho, V. Rodellar, A. Ál-
1313 varez, R. Martínez-Olalla, Evaluation of Voice Pathology Based on the Estimation of Vocal Fold
1314 Biomechanical Parameters, *Journal of Voice* 21 (4) (2007) 450–476.
- 1315 [173] C. Carmona-Duarte, R. Plamondon, P. Gómez-Vilda, M. A. Ferrer, J. B. Alonso, A. R. M. Londral,
1316 Application of the Lognormal Model to the Vocal Tract Movement to Detect Neurological Diseases in
1317 Voice, in: *Smart Innovation, Systems and Technologies, Vol. 60, 2016*, pp. 25–35.
- 1318 [174] P. Gómez-Vilda, R. Fernández-Baillo, V. Rodellar-Biarge, V. N. Lluís, A. Álvarez-Marquina, L. M.
1319 Mazaira-Fernández, R. Martínez-Olalla, J. I. Godino-Llorente, Glottal Source biometrical signature
1320 for voice pathology detection, *Speech Communication* 51 (9) (2009) 759–781.

- 1321 [175] J. Godino-Llorente, S. Aguilera-Navarro, P. Gomez-Vilda, Automatic detection of voice impairments
1322 due to vocal misuse by means of Gaussian mixture models, in: 2001 Conference Proceedings of the
1323 23rd Annual International Conference of the IEEE Engineering in Medicine and Biology Society, IEEE,
1324 2001, pp. 1723–1726.
- 1325 [176] R. F. Leonarduzzi, G. a. Alzamendi, G. Schlotthauer, M. E. Torres, Wavelet leader multifractal analysis
1326 of period and amplitude sequences from sustained vowels, *Speech Communication* 72 (2015) 1–12.
- 1327 [177] D. D. Mehta, M. Zañartu, S. W. Feng, H. A. Cheyne, R. E. Hillman, Mobile Voice Health Monitoring
1328 Using a Wearable Accelerometer Sensor and a Smartphone Platform, *IEEE Transactions on Biomedical*
1329 *Engineering* 59 (11) (2012) 3090–3096.
- 1330 [178] P. Aichinger, I. Roesner, B. Schneider-Stickler, M. Leonhard, D.-M. Denk-Linnert, W. Bigenzahn, A. K.
1331 Fuchs, M. Hagmüller, G. Kubin, Towards Objective Voice Assessment: The Diplophonia Diagram,
1332 *Journal of Voice* 31 (2) (2017) 253.e17–253.e26.
- 1333 [179] M. Döllinger, M. Kunduk, M. Kaltenbacher, S. Vondenhoff, A. Ziethe, U. Eysholdt, C. Bohr, Analysis
1334 of vocal fold function from acoustic data simultaneously recorded with high-speed endoscopy, *Journal*
1335 *of Voice* 26 (6) (2012) 726–733.
- 1336 [180] L. Verde, G. De Pietro, G. Sannino, A methodology for voice classification based on the personalized
1337 fundamental frequency estimation, *Biomedical Signal Processing and Control* 42 (2018) 134–144.
- 1338 [181] a. M. Sulter, H. P. Wit, H. K. Schutte, D. G. Miller, A structured approach to voice range profile
1339 (phonetogram) analysis., *Journal of speech and hearing research* 37 (5) (1994) 1076–85.
- 1340 [182] A. Behrman, C. J. Agresti, E. Blumstein, G. Sharma, Meaningful features of voice range profiles
1341 from patients with organic vocal fold pathology: A preliminary study, *Journal of Voice* 10 (3) (1996)
1342 269–283.
- 1343 [183] A. E. Hallin, K. Fröst, E. B. Holmberg, M. Södersten, Voice and speech range profiles and Voice
1344 Handicap Index for males – methodological issues and data, *Logopedics Phoniatrics Vocology* 37 (2)
1345 (2012) 47–61.
- 1346 [184] J. Goddard, G. Schlotthauer, M. Torres, H. Rufiner, Dimensionality reduction for visualization of
1347 normal and pathological speech data, *Biomedical Signal Processing and Control* 4 (3) (2009) 194–201.
- 1348 [185] M. FarruĀs, J. Hernando, Using Jitter and Shimmer in speaker verification, *IET Signal Processing*
1349 3 (November 2008) (2009) 247.
- 1350 [186] M. Shu, J. J. Jiang, M. Willey, The Effect of Moving Window on Acoustic Analysis, *Journal of Voice*
1351 30 (1) (2016) 5–10.

- 1352 [187] V. Uloza, E. Padervinskis, A. Vegiene, R. Pribuisiene, V. Saferis, E. Vaiciukynas, A. Gelzinis,
1353 A. Verikas, Exploring the feasibility of smart phone microphone for measurement of acoustic voice
1354 parameters and voice pathology screening, *European Archives of Oto-Rhino-Laryngology* 272 (11)
1355 (2015) 3391–3399.
- 1356 [188] Y. Zhang, J. J. Jiang, Acoustic Analyses of Sustained and Running Voices From Patients With Laryn-
1357 geal Pathologies, *Journal of Voice* 22 (1) (2008) 1–9.
- 1358 [189] S. Bielamowicz, J. Kreiman, B. R. Gerratt, M. S. Dauer, G. S. Berke, Comparison of voice analysis
1359 systems for perturbation measurement, *Journal of Speech & Hearing Research* 39 (1) (1996) 126–134.
- 1360 [190] P. R. Scalassara, C. D. Maciel, R. C. Guido, J. C. Pereira, E. S. Fonseca, A. N. Montagnoli, S. B.
1361 Júnior, L. S. Vieira, F. L. Sanchez, Autoregressive decomposition and pole tracking applied to vocal
1362 fold nodule signals, *Pattern Recognition Letters* 28 (11) (2007) 1360–1367.
- 1363 [191] A. Tsanas, M. a. Little, P. E. McSharry, J. Spielman, L. O. Ramig, Novel speech signal processing
1364 algorithms for high-accuracy classification of Parkinsons disease, *IEEE Transactions on Biomedical*
1365 *Engineering* 59 (5) (2012) 1264–1271.
- 1366 [192] A. Tsanas, M. A. Little, P. E. McSharry, L. O. Ramig, Nonlinear speech analysis algorithms mapped
1367 to a standard metric achieve clinically useful quantification of average Parkinson’s disease symptom
1368 severity, *Journal of The Royal Society Interface* 8 (59) (2011) 842–855.
- 1369 [193] E. Yumoto, W. J. Gould, T. Baer, Harmonics to Noise Ratio as hoarseness index of degree of hoarseness,
1370 *Journal of the Acoustical Society of America* 71 (6) (1982) 1544–1549.
- 1371 [194] D. Michaelis, T. Gramss, H. W. Strube, Glottal-to-noise excitation ratio—a new measure for de-
1372 scribing pathological voices, *Acta Acustica united with Acustica* 83 (4) (1997) 700–706.
- 1373 [195] G. de Krom, A cepstrum-based technique for determining a harmonics-to-noise ratio in speech signals.,
1374 *Journal of speech, language, and hearing research* 36 (2) (1993) 254–266.
- 1375 [196] H. Kasuya, Normalized noise energy as an acoustic measure to evaluate pathologic voice, *The Journal*
1376 *of the Acoustical Society of America* 80 (5) (1986) 1329.
- 1377 [197] C. Peng, W. Chen, X. Zhu, B. Wan, D. Wei, B. Engineering, Pathological Voice Classification Based
1378 on a Single Vowel ’ s Acoustic Features, *IEEE Seventh International Conference on Computer and*
1379 *Information Technology* (2007) 1106–1110.
- 1380 [198] W. S. Winholtz, L. O. Ramig, Vocal tremor analysis with the Vocal Demodulator., *Journal of speech*
1381 *and hearing research* 35 (3) (1992) 562–73.
- 1382 [199] A. Akbari, M. K. Arjmandi, Employing linear prediction residual signal of wavelet sub-bands in auto-
1383 matic detection of laryngeal pathology, *Biomedical Signal Processing and Control* 18 (2015) 293–302.

- 1384 [200] Z. Ali, G. Muhammad, M. F. Alhamid, An Automatic Health Monitoring System for Patients Suffering
1385 from Voice Complications in Smart Cities, *IEEE Access* 5.
- 1386 [201] M. Alsulaiman, Voice pathology assessment systems for dysphonic patients: Detection, classification,
1387 and speech recognition, *IETE Journal of Research* 60 (2) (2014) 156–167.
- 1388 [202] J. I. Godino-Llorente, S. Aguilera-Navarro, P. Gómez-Vilda, PC, LPCC and MFCC Parameterisation
1389 Applied to the Detection of Voice Impairments, in: *Sixth International Conference on Spoken Language*
1390 *Processing*, 2000.
- 1391 [203] J. C. Saldanha, T. Ananthakrishna, R. Pinto, Vocal Fold Pathology Assessment Using Mel-Frequency
1392 Cepstral Coefficients and Linear Predictive Cepstral Coefficients Features, *Journal of Medical Imaging*
1393 *and Health Informatics* 4 (2) (2014) 168–173.
- 1394 [204] Z. Ali, M. Alsulaiman, G. Muhammad, I. Elamvazuthi, A. Al-nasheri, T. A. Mesallam, M. Farahat,
1395 K. H. Malki, Intra- and Inter-database Study for Arabic, English, and German Databases: Do Con-
1396 ventional Speech Features Detect Voice Pathology?, *Journal of Voice* 31 (3) (2017) 386.e1–386.e8.
- 1397 [205] M. S. Hossain, G. Muhammad, A. Alamri, Smart healthcare monitoring: a voice pathology detection
1398 paradigm for smart cities, *Multimedia Systems* (2017) 1–11.
- 1399 [206] F. Amara, M. Fezari, H. Bourouba, An Improved GMM-SVM System based on Distance Metric for
1400 Voice Pathology Detection, *Appl. Math. Inf. Sci* 10 (3) (2016) 1061–1070.
- 1401 [207] H. Florian, G. Stemmer, C. Hacker, F. Brugnara, Revising Perceptual Linear Prediction (PLP),
1402 *Interspeech 2005* (2005) 2997–3000.
- 1403 [208] H. Hermansky, Perceptual linear predictive (PLP) analysis of speech, *The Journal of the Acoustical*
1404 *Society of America* 87 (1990) 1738.
- 1405 [209] A. Benba, A. Jilbab, A. Hammouch, Discriminating Between Patients With Parkinson’s and Neuro-
1406 logical Diseases Using Cepstral Analysis, *IEEE Transactions on Neural Systems and Rehabilitation*
1407 *Engineering* 24 (10) (2016) 1100–1108.
- 1408 [210] J. R. Orozco-Aroyave, F. Höning, J. D. Arias-Londoño, J. F. Vargas-Bonilla, E. Nöth, Spectral and
1409 cepstral analyses for Parkinson’s disease detection in Spanish vowels and words, *Expert Systems* 32 (6)
1410 (2015) 688–697.
- 1411 [211] A. Akbari, M. K. Arjmandi, An efficient voice pathology classification scheme based on applying multi-
1412 layer linear discriminant analysis to wavelet packet-based features, *Biomedical Signal Processing and*
1413 *Control* 10 (1) (2014) 209–223.
- 1414 [212] H. Khadivi Heris, B. Seyed Aghazadeh, M. Nikkhah-Bahrami, Optimal feature selection for the assess-
1415 ment of vocal fold disorders, *Computers in Biology and Medicine* 39 (10) (2009) 860–868.

- 1416 [213] K. Umaphathy, S. Krishnan, Feature analysis of pathological speech signals using local discriminant
1417 bases technique., *Medical & biological engineering & computing* 43 (4) (2005) 457–64.
- 1418 [214] K. Umaphathy, S. Krishnan, V. Parsa, D. G. Jamieson, Discrimination of pathological voices using a
1419 time-frequency approach., *IEEE transactions on bio-medical engineering* 52 (3) (2005) 421–30.
- 1420 [215] N. Erfanian Saeedi, F. Almasganj, F. Torabinejad, Support vector wavelet adaptation for pathological
1421 voice assessment, *Computers in Biology and Medicine* 41 (9) (2011) 822–828.
- 1422 [216] B. Ghoraani, K. Umaphathy, L. Sugavaneswaran, S. Krishnan, Pathological Speech Signal Analysis
1423 Using Time-Frequency Approaches, *Critical Reviews in Biomedical Engineering* 40 (1) (2012) 63–95.
- 1424 [217] J. Hillenbrand, R. A. Cleveland, R. L. Erickson, Acoustic Correlates of Breathly Vocal Quality, *Journal*
1425 *of Speech Language and Hearing Research* 37 (4) (1994) 769–778.
- 1426 [218] J. Hillenbrand, R. A. Houde, Acoustic Correlates of Breathly Vocal Quality: Dysphonic Voices and
1427 Continuous Speech, *Journal of Speech Language and Hearing Research* 39 (2) (1996) 311.
- 1428 [219] P. Grassberger, Randomness, Information, and Complexity, *Complexity* (1) (2012) 1–20.
- 1429 [220] M. Costa, A. Goldberger, C. Peng, Multiscale entropy analysis of biological signals, *Physical Review*
1430 *E* 71 (2) (2005) 1–18.
- 1431 [221] M. D. Costa, C.-K. Peng, A. L. Goldberger, Multiscale Analysis of Heart Rate Dynamics: Entropy
1432 and Time Irreversibility Measures, *Cardiovascular Engineering* 8 (2) (2008) 88–93.
- 1433 [222] I. R. Titze, *The Myoelastic Aerodynamic Theory of Phonation*, 1st Edition, National Center for Voice
1434 and Speech, 2006.
- 1435 [223] I. Steinecke, H. Herzel, Bifurcations in an asymmetric vocal-fold model., *The Journal of the Acoustical*
1436 *Society of America* 97 (3) (1995) 1874–84.
- 1437 [224] M. A. Little, P. E. McSharry, S. J. Roberts, D. A. Costello, I. M. Moroz, Exploiting Nonlinear Re-
1438 currence and Fractal Scaling Properties for Voice Disorder Detection, *BioMedical Engineering OnLine*
1439 6 (1) (2007) 23.
- 1440 [225] R. T. Sataloff, *Laryngology (Sataloff’s Comprehensive Textbook of Otolaryngology: Head & Neck*
1441 *Surgery)*, Jaypee Brothers, Medical Publishers Pvt. Limited, 2015.
- 1442 [226] G. Vaziri, F. Almasganj, R. Behroozmand, Pathological assessment of patients’ speech signals using
1443 nonlinear dynamical analysis, *Computers in Biology and Medicine* 40 (1) (2010) 54–63.
- 1444 [227] L. Matassini, R. Hegger, H. Kantz, C. Manfredi, Analysis of vocal disorders in a feature space, *Medical*
1445 *Engineering & Physics* 22 (6) (2000) 413–418.

- 1446 [228] J. Arias-Londoño, J. Godino-Llorente, G. Castellanos-Domínguez, N. Sáenz-Lechón, V. Osma-
1447 Ruiz, Complexity analysis of pathological voices by means of hidden markov entropy measurements,
1448 in: 2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society,
1449 IEEE, 2009, pp. 2248–2251.
- 1450 [229] Z. Ali, I. Elamvazuthi, M. Alsulaiman, G. Muhammad, Detection of Voice Pathology using Fractal
1451 Dimension in a Multiresolution Analysis of Normal and Disordered Speech Signals, *Journal of medical*
1452 *systems* 40 (1) (2016) 20.
- 1453 [230] J. A. Gómez-García, J. I. Godino-Llorente, G. Castellanos-Domínguez, Non uniform Embedding based
1454 on Relevance Analysis with reduced computational complexity: Application to the detection of patholo-
1455 gies from biosignal recordings, *Neurocomputing* 132 (0) (2014) 148–158.
- 1456 [231] J.-R. Orozco Arroyave, J.-F. Vargas Bonilla, E. Delgado Trejos, Acoustic Analysis and Non Linear
1457 Dynamics Applied to Voice Pathology Detection: A Review, *Recent Patents on Signal Processing*
1458 2 (2) (2012) 96–107.
- 1459 [232] L. Moro-Velázquez, J. A. Gómez-García, J. I. Godino-Llorente, Voice Pathology Detection Using Mod-
1460 ulation Spectrum-Optimized Metrics, *Frontiers in Bioengineering and Biotechnology* 4 (1).
- 1461 [233] Z. Ali, M. Talha, M. Alsulaiman, A Practical Approach: Design and Implementation of a Healthcare
1462 Software for Screening of Dysphonic Patients, *IEEE Access* 3536 (c) (2017) 1–1.
- 1463 [234] Y. Maryn, M. De Bodt, N. Roy, The Acoustic Voice Quality Index: Toward improved treatment
1464 outcomes assessment in voice disorders, *Journal of Communication Disorders* 43 (3) (2010) 161–174.
- 1465 [235] M. Fröhlich, D. Michaelis, H. W. Strube, E. Kruse, Acoustic voice quality description : Case studies for
1466 different regions of the hoarseness diagram, *Advances in Quantitative Laryngoscopy*, 2nd Round
1467 Table (1997) 143–150.
- 1468 [236] A. Fontes, P. Souza, A. Neto, A. Martins, L. Silveira, Classification System of Pathological Voices
1469 Using Correntropy, *Mathematical Problems in Engineering* 2014 (2014) 1–7.
- 1470 [237] P. Gómez-Vilda, E. S. Segundo, L. M. Mazaira-Fernández, A. Álvarez-Marquina, V. Rodellar-Biarge,
1471 Using Dysphonic Voice to Characterize Speaker’s Biometry, *Language and Law / Linguagem e Direito*
1472 1 (2) (2014) 42–66.
- 1473 [238] D. Hemmerling, A. Skalski, J. Gajda, Voice data mining for laryngeal pathology assessment, *Computers*
1474 *in Biology and Medicine* 69 (2016) 270–276.
- 1475 [239] M. K. Arjmandi, M. Pooyan, M. Mikaili, M. Vali, A. Moqarehzadeh, Identification of voice disorders
1476 using long-time features and support vector machine with different feature reduction methods, *Journal*
1477 *of Voice* 25 (6) (2011) e275–e289.

- 1478 [240] Z. Ali, M. Alsulaiman, I. Elamvazuthi, G. Muhammad, T. A. Mesallam, M. Farahat, K. H. Malki, Voice
1479 pathology detection based on the modified voice contour and SVM, *Biologically Inspired Cognitive*
1480 *Architectures* 15 (2016) 10–18.
- 1481 [241] S. Shilaskar, A. Ghatol, P. Chatur, Medical decision support system for extremely imbalanced datasets,
1482 *Information Sciences* 384 (2017) 205–219.
- 1483 [242] N. Erfanian Saeedi, F. Almasganj, Wavelet adaptation for automatic voice disorders sorting, *Computers*
1484 *in Biology and Medicine* 43 (6) (2013) 699–704.
- 1485 [243] M. Markaki, Y. Stylianou, Normalized modulation spectral features for cross-database voice pathol-
1486 ogy detection, in: *Proceedings of the Annual Conference of the International Speech Communication*
1487 *Association, INTERSPEECH, 2009*, pp. 935–938.
- 1488 [244] T. Fawcett, An introduction to ROC analysis, *Pattern Recognition Letters* 27 (8) (2006) 861–874.
- 1489 [245] T. Dunstone, N. Yager (Eds.), *Biometric System and Data Analysis*, Springer US, Boston, MA, 2009.
- 1490 [246] A. P. Bradley, The use of the area under the ROC curve in the evaluation of machine learning algo-
1491 rithms, *Pattern Recognition* 30 (7) (1997) 1145–1159.
- 1492 [247] R. Fernández-Pozo, J. L. Blanco-Murillo, L. Hernández Gómez, E. López, J. Alcázar-Ramírez, D. Torre-
1493 Toledano, Assessment of Severe Apnoea through Voice Analysis, *Automatic Speech, and Speaker*
1494 *Recognition Techniques, EURASIP Journal on Advances in Signal Processing* 2009 (1) (2009) 1–12.
- 1495 [248] A. Montero Benavides, J. L. Blanco Murillo, R. Fernández Pozo, F. Espinoza Cuadros, D. Torre
1496 Toledano, J. D. Alcázar-Ramírez, L. A. Hernández Gómez, Formant Frequencies and Bandwidths in
1497 Relation to Clinical Variables in an Obstructive Sleep Apnea Population, *Journal of Voice* 30 (1) (2016)
1498 21–29.
- 1499 [249] J. Alcázar-Ramírez, R. Fernández-Pozo, J. L. Blanco-Murillo, L. Hernández, L. López, F. Linde,
1500 D. Torre-Toledano, *Automatic Speaker Recognition Techniques: A new Tool for Sleep Apnoea Diag-*
1501 *nosis, Am. J. Respir. Crit. Care Med.*
- 1502 [250] J. R. Orozco-Arroyave, J. F. Vargas-Bonilla, J. D. Arias-Londoño, S. Murillo-Rendón, G. Castellanos-
1503 Domínguez, J. F. Garcés, Nonlinear Dynamics for Hypernasality Detection in Spanish Vowels and
1504 Words, *Cognitive Computation* 5 (4) (2013) 448–457.
- 1505 [251] P. Gómez-Vilda, V. Rodellar-Biarge, V. Nieto-Lluis, K. L. de Ipiña, A. Álvarez-Marquina, R. Martínez-
1506 Olalla, M. Ecay-Torres, P. Martínez-Lage, Phonation biomechanic analysis of Alzheimer’s Disease
1507 cases, *Neurocomputing* 167 (2015) 83–93.

- 1508 [252] K. López-de Ipiña, J.-B. Alonso, C. Travieso, J. Solé-Casals, H. Egiraun, M. Faundez-Zanuy, A. Ezeiza,
1509 N. Barroso, M. Ecay-Torres, P. Martinez-Lage, U. Lizardui, On the Selection of Non-Invasive Methods
1510 Based on Speech Analysis Oriented to Automatic Alzheimer Disease Diagnosis, *Sensors* 13 (5) (2013)
1511 6730–6745.
- 1512 [253] J. Ortega-Garcia, J. Gonzalez-Rodriguez, V. Marrero-Aguiar, AHUMADA: A large speech corpus in
1513 spanish for speaker characterization and identification, *Speech Communication* 31 (2-3) (2000) 255–264.
- 1514 [254] S. L. Velleman, B. Z. Pearson, Differentiating Speech Sound Disorders From Phonological Dialect
1515 Differences: Implications for Assessment and Intervention, *Top Language Disorders* 30 (3) (2010)
1516 176–188.
- 1517 [255] C. Huang, T. Chen, S. Z. Li, E. Chang, J.-L. Zhou, Analysis of speaker variability., in: *Proc. INTER-*
1518 *SPEECH*, no. 49, 2001, pp. 1377–1380.
- 1519 [256] C. Huang, T. Chen, E. Chang, Accent Issues in Large Vocabulary Continuous Speech Recognition,
1520 *International Journal of Speech Technology* 7 (2/3) (2004) 141–153.
- 1521 [257] M. Benzeghiba, R. De Mori, O. Deroo, S. Dupont, T. Erbes, D. Jovet, L. Fissore, P. Laface, a. Mertins,
1522 C. Ris, R. Rose, V. Tyagi, C. Wellekens, Automatic speech recognition and speech variability: A review,
1523 *Speech Communication* 49 (10-11) (2007) 763–786.
- 1524 [258] H. Traunmüller, A. Eriksson, Acoustic effects of variation in vocal effort by men, women, and children,
1525 *The Journal of the Acoustical Society of America* 107 (6) (2000) 3438.
- 1526 [259] A. L. Rosenthal, S. Y. Lowell, R. H. Colton, Aerodynamic and Acoustic Features of Vocal Effort,
1527 *Journal of Voice* 28 (2) (2014) 144–153.
- 1528 [260] S. N. Awan, A. Giovinco, J. Owens, Effects of Vocal Intensity and Vowel Type on Cepstral Analysis
1529 of Voice, *Journal of Voice* 26 (5) (2012) 670.e15–670.e20.
- 1530 [261] D. Z. Huang, F. D. Minifie, H. Kasuya, S. X. Lin, Measures of vocal function during changes in vocal
1531 effort level, *Journal of Voice* 9 (4) (1995) 429–438.
- 1532 [262] M. Brockmann, C. Storck, P. N. Carding, M. J. Drinnan, Voice Loudness and Gender Effects on Jitter
1533 and Shimmer in Healthy Adults, *Journal of Speech, Language, and Hearing Research* 51 (5) (2008)
1534 1152–1160.
- 1535 [263] M. Brockmann, M. J. Drinnan, C. Storck, P. N. Carding, Reliable Jitter and Shimmer Measurements
1536 in Voice Clinics: The Relevance of Vowel, Gender, Vocal Intensity, and Fundamental Frequency Effects
1537 in a Typical Clinical Task, *Journal of Voice* 25 (1) (2011) 44–53.
- 1538 [264] P. H. Dejonckere, Effect of louder voicing on acoustical measurements in dysphonic patients, *Logopedics*
1539 *Phoniatrics Vocology* 23 (2) (1998) 79–84.

- 1540 [265] L. Chen, X. Mao, Y. Xue, L. L. Cheng, Speech emotion recognition: Features and classification models,
1541 Digital Signal Processing 22 (6) (2012) 1154–1160.
- 1542 [266] S. Steidl, A. Batliner, D. Seppi, B. Schuller, On the Impact of Children’s Emotional Speech on Acoustic
1543 and Language Models, EURASIP Journal on Audio, Speech, and Music Processing 2010 (14) (2010)
1544 1–14.
- 1545 [267] B. Schuller, J. Stadermann, G. Rigoll, in: Speech Prosody, 2006.
- 1546 [268] D. Childers, K. Wu, Automatic recognition of gender by voice, in: International Conference on Acous-
1547 tics, Speech, and Signal Processing (ICASSP), 1988, pp. 603–606.
- 1548 [269] W. D. Andrews, M. A. Kohler, J. P. Campbell, J. J. Godfrey, J. Hernández-Cordero, Gender-dependent
1549 phonetic refraction for speaker recognition, in: International Conference on Acoustics, Speech, and
1550 Signal Processing (ICASSP), Vol. 1, IEEE, 2002, pp. I–149.
- 1551 [270] T. J. Hixon, G. Weismer, J. D. Hoit, Preclinical Speech Science: Anatomy, Physiology, Acoustics, and
1552 Perception, 2nd Edition, Plural Publishing, 2013.
- 1553 [271] A. Behrman, Speech and Voice Science, Plural Publishing Inc, 2007.
- 1554 [272] M. Södersten, S. Hertegård, B. Hammarberg, Glottal closure, transglottal airflow, and voice quality
1555 in healthy middle-aged women., Journal of Voice 9 (2) (1995) 182–197.
- 1556 [273] H. M. Hanson, E. S. Chuang, Glottal characteristics of male speakers: acoustic correlates and compar-
1557 ison with female data., The Journal of the Acoustical Society of America 106 (2) (1999) 1064–77.
- 1558 [274] D. H. Klatt, L. C. Klatt, Analysis, synthesis, and perception of voice quality variations among female
1559 and male talkers, The Journal of the Acoustical Society of America 87 (2) (1990) 820–857.
- 1560 [275] K. Wu, D. G. Childers, Gender recognition from speech. Part I: Coarse analysis., The Journal of the
1561 Acoustical Society of America 90 (4 Pt 1) (1991) 1828–40.
- 1562 [276] E. Mendoza, N. Valencia, J. Muñoz, H. Trujillo, Differences in voice quality between men and women:
1563 use of the long-term average spectrum (LTAS)., Journal of Voice 10 (1) (1996) 59–66.
- 1564 [277] H. M. Hanson, Glottal characteristics of female speakers: acoustic correlates., The Journal of the
1565 Acoustical Society of America 101 (1) (1997) 466–81.
- 1566 [278] E. Holmberg, R. Hillman, J. Perkell, Glottal airflow and transglottal air pressure measurements for
1567 male and female speakers in soft, normal, and loud voice, The Journal of the Acoustical Society of
1568 America 84 (2) (1988) 511–529.
- 1569 [279] R. B. Monsen, a. M. Engebretson, Study of variations in the male and female glottal wave., The Journal
1570 of the Acoustical Society of America 62 (4) (1977) 981–93.

- 1571 [280] J. Abitbol, P. Abitbol, B. Abitbol, Sex hormones and the female voice., *Journal of Voice* 13 (3) (1999)
1572 424–46.
- 1573 [281] B. Gold, N. Morgan, D. Ellis, *Speech and Audio Signal Processing: Processing and Perception of*
1574 *Speech and Music*, Wiley-Interscience, 2011.
- 1575 [282] D. D. Deliyski, H. S. Shaw, M. K. Evans, Adverse Effects of Environmental Noise on Acoustic Voice
1576 Quality Measurements, *Journal of Voice* 19 (1) (2005) 15–28. doi:10.1016/j.jvoice.2004.07.003.
- 1577 [283] A. H. Poorjam, J. R. Jensen, M. A. Little, M. G. Christensen, Dominant Distortion Classification for
1578 Pre-Processing of Vowels in Remote Biomedical Voice Analysis, in: *Interspeech 2017*, Vol. 2017-Augus,
1579 2017, pp. 289–293.