Rita Singh

# Profiling Humans from their Voice

# Profiling Humans from their Voice

Rita Singh

# Profiling Humans from their Voice

Springer

Rita Singh
Carnegie Mellon University
Pittsburgh, PA, USA

"I HAVE SHOWN THAT VOICE QUALITY
DEPENDS ABSOLUTELY UPON BONE
STRUCTURE, THAT BONE STRUCTURE IS
INHERITED, AND THAT THEREFORE VOCAL
QUALITY IS INHERITED."

—Walter B. Swift, 1916

*The possibility of voice inheritance. In Review of Neurology and
Psychiatry, Volume 14, Page 106.*

# Preface

## What is Different

Computational profiling as described in this book differs from prior efforts in one key respect. The fact that voice is correlated to many human parameters has been known to scientists, doctors, philosophers, priests, performers, soothsayers—in fact, to people of myriad vocations for centuries. Many of these correlations have been scientifically investigated, and demonstrated or proven in many ways. The key difference between that earlier science, and what this book represents is simply this: while earlier methods were focused on demonstration and proof of existence (of such relationships), and used observable relations to attempt to predict human parameters from voice, the current science does not require such observables. It is built instead on the hypothesis that *if any factor whatsoever influences the human mind or body, and if that influence can be linked to the human voice production mechanism through any pathway whatsoever, then there **must** exist an effect on voice*. The current science of profiling is then all about discovering those effects. If we hypothesize the existence of some influence, whose effects we cannot model or observe through standard mechanisms known to us, then we must devise artificially intelligent mechanisms that *can* model or observe those. This basically represents a handover of the capability of discovery to intelligent systems designed by us.

This book traverses pathways hewn through information in multiple disciplines, and represents one journey undertaken in search of solutions to such mysteries of the human voice.

## About This Book

The ability to shape sound is an amazing gift that most intelligent creatures have. Some do this in response to stimuli, others do this to convey meaning as well. With the evolution of intellectual abilities, it was only natural for this ability to evolve to

convey more complex thoughts and deeper meaning. It was also natural, through the process of natural selection, for those kinds of information-embedding mechanisms that were more supportive of survival to remain, propagate, and be refined. Sound is no exception. Deliberately or inadvertently, partly for evolutionary reasons and partly due to our specific biological construction, information about each person is embedded in their voice.

At this point in time, we don't have a good idea of just how much information is embedded in the human voice. Reasoning about speech as a biomechanical, social, and cognitive process leads us to believe that there is a tremendous amount of information in it—more than we are capable of assimilating through our limited capabilities of auditory observation and perception.

This book captures some of my thoughts, ideas, and research on discovering or even guessing the range and extent of information embedded in the human voice, on deriving it quantitatively from voice signals, and using it to infer bio-relevant facts about the speaker and their environment. In my experience, this endeavor is so rife with challenges, and human voice is of such tremendous importance, that it deserves to be assigned the status of a subfield of acoustic intelligence in its own right—so I call it *Profiling Humans from their Voice*, which is also the title of this book.

The book has two parts. Part I takes a sweeping look at the landscape of scientific explorations into the human voice, which has been the subject of an astounding volume of research and observation, literally over centuries. Voice, its acoustics, its content, the effect of various factors on these, and conversely *its* effect on them, and on other humans, its perceptions, and its manipulations are all discussed in this part. This part also dives into the voice *signal* from a signal processing perspective. It (very) briefly elucidates the concepts that might be relevant or foundational to profiling, attempting to link some subjective observations of the quality of voice, based on which most human judgments are made, to explainable or quantifiable signal characteristics.

Part II deals with computational profiling: the computationalization of human judgment (and beyond). Predicated largely on concepts in machine learning and artificial intelligence, this part discusses mechanisms for information discovery, feature engineering, and the deduction of profile parameters from them. It discusses the subject of reconstructing the human persona from voice, and its reverse—the reconstruction of voice from information about the human persona. It ends with a discussion of the applications and future outlook for the science of profiling, and also ethical issues associated with its use. Issues of reliability, confidence estimation, statistical verification, etc., that are extremely important in practical applications, are all discussed.

There are multiple chapters under each part, divided up to make the overall theme of the part easier to navigate. Each part is provided with a summary of the chapters it includes.

Although meant to be a technical exposition, in my mind this book is a canvas with dabs of paint from many fields. Together they form a complete picture—but it is still an underdrawing. There is just so much that one can write in a book of

normal dimensions. As I type the last full stop in this book, it feels incomplete in many respects. I hope the studies referenced in this book can fill in the missing details to some extent.

More than anything, I hope this book is both enjoyable and informative to all readers.

Pittsburgh, PA, USA                                                                  Rita Singh
March 2019

# Acknowledgements

# Contents

# List of Figures

# List of Tables

# Part I
# Profiling and the Human Voice

Part I of this book connects findings from scientific explorations in different fields to define and outline an emergent area of research—that of profiling humans from their voice. It is not a comprehensive review of these fields, but rather a cursory attempt to build a perspective that may be central to future developments in the area. There are six chapters in this part. Chapter 1 defines profiling and segregates its different aspects. The rest of the chapters provide key insights into each aspect, drawn from existing literature that often spans multiple fields. Chapter 2 is a brief overview of the key concepts in voice production and audition that have a bearing on profiling. Chapter 3 describes the relationship of voice to different categories of profile parameters. Chapters 4 and 5 deal with the information content in the voice signal, much of which is based on core signal processing concepts. These are explained where necessary within the chapters. Chapter 6 deals with voice quality, a very subjective and descriptive entity that has been so extensively correlated with human parameters that it warrants a separate chapter in itself.

# Chapter 1
# Profiling and Its Facets

## 1.1 Profiling

The term *profiling from voice* refers to the deduction of personal characteristics, and information about the circumstances and environment of a speaker from their voice. At the outset, we note the distinction between the terms *voice* and *speech*. "Voice" refers to sound produced in the human vocal tract. "Speech" is the signal produced by modulating voice into meaningful patterns.

We will use the term "*profiling*" to refer to the process of deduction of these entities, and the term "*parameters*" to refer to the entities deduced. Later in this book we will often use the term "*features*," which will exclusively refer to mathematical representations derived from the voice signal for the purpose of profiling. Other terms such as "*characteristic*," "*entity*" etc. will be used in their normal linguistic sense, depending on context. For example, properties of the voice signal, such as voice quality, speaking rate, loudness etc. will be referred to as "*signal characteristics*." In later chapters, we will also see that features are a subset of signal characteristics, but not vice-versa. As long as we clearly adhere to this specific usage of terms, the exposition in this book can be followed with precision.

We begin with expanding the definition of some of these terms, to establish the specific organization of topics and the vocabulary that we will use in this book.

### *1.1.1 Parameters*

Profiling comprises the deduction of *all* kinds of parameters from voice. This space is vast, and an initial categorization of parameters based on some broad-level commonalities is needed in order to proceed. Although it is obvious that many categorizations are possible, in this book we will follow the specific broad-level subcategorization of parameters given in the list below. This list only includes a few samples of parameters from each category, and is not an exhaustive one.

1. **Behavioral parameters**: Dominance, leadership, public and private behavior.
2. **Demographic parameters**: Race, geographical origins, level of education.
3. **Environmental parameters**: Location of the speaker at the time of speaking, objects surrounding the person at the time of speaking, devices and communication channels used for voice capture and transmission.
4. **Medical parameters**: Presence or absence of specific diseases, medications and other substances (drugs, food, intoxicants etc.) in the body, state of physical health, state of mental health, effects of trauma, effects of medical procedures, presence or absence of physical abnormalities or disabilities.
5. **Physical parameters**: Height, weight, body-shape, facial structure.
6. **Physiological parameters**: Age, hormone levels, heart rate, blood pressure.
7. **Psychological parameters**: Personality, emotions.
8. **Sociological parameters**: Social status, income, profession.

In addition to the four key terms (profiling, parameters, features, characteristics) mentioned above, the following terms may sometimes be used in this book:

- **Factors**: A generic term that refers to elements, phenomena, processes, influencing entities etc. For example: "*Atmospheric pressure is a factor that influences vocal parameters.*" The word "*factor*" may substitute for the word "*parameter*" in some contexts that do not immediately involve deduction, e.g. "humidity is also a factor that affects voice characteristics."
- **Voice-print**: Any segment of speech, regardless of its lexical or semantic content, or completeness.
- **Bio-relevant parameters** (or *bio-parameters*): All parameter categories mentioned above except environmental parameters. Note that this is a superset of parameters that have strictly biological relevance.

## *1.1.2 Features and Signal Characteristics*

The set of features includes informative representations of the voice signal. The set of signal characteristics refers to various properties of the voice signal that have traditionally been derived and studied in the context of voice. Some of these are listed below, based on how they are measured (described in detail in later chapters). Again, this list is not exhaustive:

1. **Temporal domain features**: Properties measurable in the time domain, such as signal energy, zero crossing rate, loudness, energy, speaking rate, phonation rate etc.
2. **Spectral domain features**: Properties measurable in the frequency domain, such as pitch or fundamental frequency, harmonics, spectral flux, spectral density, spectral roll-off etc.

3. **Composite domain features**: Require measurements in both frequency and time, or are derivatives of spectral and temporal domain features, such as rhythm, melody etc.
4. **Qualitative features**: Largely human-judged properties that can be assigned to a voice signal, such as voice quality and its related components—nasality, raspiness, breathiness, roughness etc.

## 1.2  A Look at the Landscape of Voice Studies

Human voice has been widely studied in the context of, and within the purview of many scientific fields. In fact, the range of studies that have been performed in this context over the past century is staggering. Within these fields, voice has been correlated to multiple parameters.

Table 1.1 shows some examples of the fields in the context of which voice studies have been performed, with one or more sample publications listed against all except those that are likely to include voice, or are obviously sound or voice related fields. This list is by no means exhaustive. There are many more fields that have performed investigations in the context of voice and speech. Collectively, these studies reveal a wealth of information that can be potentially brought to bear on the science of profiling.

### *1.2.1  Parameters That Have Been Correlated to Voice*

Studies mentioned in Table 1.1, and others that exist abundantly in the literature, fall under two broad categories. In one, *perceptions* of various entities, such as facial appearance, body-size, benevolence, competence, dominance, dynamism, personality etc. have been correlated with voice. In the other, *measurements*, or expert ratings of multiple parameters that relate to the speaker's persona, and their speaking environment have been correlated with voice. This list overlaps with the one above, and includes a multitude of diseases such as psychiatric illnesses, physical illnesses, various genetic syndromes etc. and a very wide range of other parameters such as age, attractiveness, body size and shape, dominance, emotion, behavior, personality, intellectual capacity, deception, intoxication, drug use and abuse, exposure to chemicals and toxins, injuries, trauma, sexual bias and orientation, sexual behavior, sleep deprivation etc. In addition, the list includes myriad environmental factors such as atmospheric moisture, relative humidity, odors, pollution etc.

The parameters that have been found to influence voice, and the specific characteristics of voice that reflect them, are discussed in detail in Chap. 3.

**Table 1.1** Examples of fields of science and voice related studies

| Field [studies] | |
|---|---|
| 1 Acoustics | 2 Adenology [1] |
| 3 Algedonics [2] | 4 Anaesthesiology [3] |
| 5 Anatomy [4] | 6 Andragogy [5] |
| 7 Angiology [6, 7] | 8 Anthropobiology [8–10] |
| 9 Anthropology [11, 12] | 10 Aromachology [13] |
| 11 Archaeology [14] | 12 Astheniology [15, 16] |
| 13 Audiology | 14 Auxology [17] |
| 15 Bioecology [18] | 16 Biology [19] |
| 17 Psychiatry [20, 21] | 18 Biometrics |
| 19 Bionomics [22] | 20 Cardiology |
| 21 Catacoustics | 22 Catechectics |
| 23 Cell biology | 24 Characterology [11, 23] |
| 25 Chronobiology [24, 25] | 26 Criminology [26] |
| 27 Cinematology [27] | 28 Demography [28] |
| 29 Demology [29] | 30 Dialectology |
| 31 Dramaturgy [30, 31] | 32 Ecology [32, 33] |
| 33 Endocrinology [34, 35] | 34 Enzymology [36] |
| 35 Ephebiatrics | 36 Epidemiology [37] |
| 37 Epileptology [38] | 38 Ergonomics [39] |
| 39 Ethnology [40] | 40 Ethology [41] |
| 41 Forensics [42, 43] | 42 Gastroenterology [44] |
| 43 Genealogy | 44 Gerontology [45] |
| 45 Glossology | 46 Harmonics |
| 47 Heredity [46, 47] | 48 Histopathology [48, 49] |
| 49 Hydrodynamics [50, 51] | 50 Hydrokinetics [52, 53] |
| 51 Hydrometeorology [54] | 52 Hymnology [55] |
| 53 Hypnology [56] | 54 Immunology [57, 58] |
| 55 Immunopathology [59, 60] | 56 Kalology [61] |
| 57 Kinematics [62] | 58 Kinesics [63] |
| 59 Kinesiology [64] | 60 Koniology [65] |
| 61 Laryngology | 62 Linguistics |
| 63 Magnetics [66] | 64 Mammalogy [67] |
| 65 Mechanics [68] | 66 Metaphysics [69, 70] |
| 67 Microanatomy [71] | 68 Musicology |
| 69 Myology [72] | 70 Naology [73, 74] |
| 71 Nasology [75] | 72 Neonatology [76, 77] |
| 73 Nephrology [78] | 74 Neurobiology [79] |

**Table 1.1**  (continued)

| Field [studies] | |
|---|---|
| 75 Neurology [80, 81] | 76 Neuropsychology [82] |
| 77 Odontology [83] | 78 Olfactology [13] |
| 79 Osteology [84] | 80 Otology [85] |
| 81 Otorhinolaryngology [86, 87] | 82 Paedology |
| 83 Peridontology | 84 Parapsychology [88–91] |
| 85 Pathology | 86 Pedagogy [92, 93] |
| 87 Pharmacology [94] | 88 Pharyngology |
| 89 Phenomenology [95] | 90 Philology [96] |
| 91 Philosophy [97, 98] | 92 Phoniatrics |
| 93 Phonology | 94 Phraseology |
| 95 Physics | 96 Physiology |
| 97 Primatology [99] | 98 Prosody |
| 99 Pseudology [100] | 100 Psychobiology [101] |
| 101 Psychogenetics [102] | 102 Psychognosy [103, 104] |
| 103 Psychology | 104 Psychopathology [105] |
| 105 Psychophysics [106] | 106 Quinology [107] |
| 107 Reflexology [108] | 108 Rheumatology [109, 110] |
| 109 Rhinology [85] | 110 Semantics [111] |
| 111 Sexology [112, 113] | 112 Sociobiology [114] |
| 113 Sociology [103, 115] | 114 Somatology [116] |
| 115 Sophiology [117] | 116 Stomatology [118, 119] |
| 117 Teratology [120] | 118 Thermokinetics [121] |
| 119 Thermology [122] | 120 Threpsology [123, 124] |
| 121 Tonetics [125] | 122 Toxicology [126] |
| 123 Traumatology [127, 128] | 124 Typhlology [129–131] |
| 125 Victimology [132, 133] | 126 Virology [134, 135] |

## 1.3  Profiling Humans, by Humans

Section 1.2.1 above lists many parameters that have been correlated to voice. This is not necessarily a list of *judgments* we make from voice—a subject that has been investigated in many studies. These in turn must be differentiated from studies on what (primarily emotions and thoughts) can be *evoked* in humans by voice, although both are related in obvious ways.

### 1.3.1  Judgments Made from Voice

Even elephants can judge age, ethnicity and gender of humans from their voices [136].

It is well known that we make myriad judgments about people based on voice. In the field of speech perception, this aspect has been especially well studied. Interestingly, the ability to make judgments about people from their voice differs between the genders [137].

Of the judgments made, many are made independently of the co-occurrence of other judgments (such as height and education, which have no correlation with one another), while some are conditioned on other parameters. For example, the judgment of gender is conditioned on age (it is difficult to tell the gender of very young children from their voice). The properties of the voice signal that play the greatest role in allowing us to make these judgments are those related to *voice quality* [138]. Voice quality is not a single entity—rather it is a loosely defined term given to a collection of entities that comprise it. We discuss voice quality and its measurement in detail in later chapters of this book. The list below shows a number of judgments that are made based on the perception of voice *quality* alone, as in [139]. Note that the categorizations here differ from what we will use in later chapters.

- **Physical characteristics**: age, appearance (height, weight, attractiveness), dental and oral/nasal status, health status, fatigue, identity, intoxication, race, ethnicity, gender, sexual orientation, smoking habits.
- **Psychological characteristics**: arousal (relaxed/hurried), competence, emotional status/mood, intelligence, personality, psychiatric status, stress, truthfulness.
- **Social characteristics**: education, occupation, regional origin, role in conversational setting, social status.

Because humans ubiquitously and subconsciously judge personality from voices, people with voice disorders often fear that their personalities may be incorrectly portrayed by their voice, and refrain from speaking in many instances. This is a well-known problem encountered in people with voice disorders. Of course, human hearing is not perfect, and many of the impressions we form of other people's characteristics based on even normal voices are inaccurate. This is also something that we experience quite often in today's world—when we talk to strangers over phone and meet them later only to discover that they don't fit all of the impressions we formed of them from their voice.

Humans also make significant judgments about the *environment* of the speaker from their voice. For example, we can easily tell whether a recording has been transmitted over the phone, or recorded in close (physical) vicinity of the speaker. We "hear" echoes and reverberations in the acoustic signal, and its effect on voice, and are able to gauge the size of the enclosure within which it was produced to coarse degrees from voice (small space, large hall etc.). Often when we hear the voice of a person calling from a moving vehicle (e.g. car, train, boat), from the manner in which the voice "shakes" (in addition to other acoustic cues) we are able to tell that the person is on a moving vehicle, and are often able to identify the broad category of the vehicle. Many of the judgments that people make about their environment are not well documented or studied. Nevertheless, it is clear that the voice signal carries cues about the environment as well.

Human abilities to judge from voice are thought to arise from a long evolutionary process, wherein they were germane to survival. In animals, such judgment co-evolved with the successful use of vocal cues to convey different *impressions*, such as those of danger, interest, dominance. Studies have found that many animal species use vocal cues to indicate situations that are threatening, or relevant to mating and social relationships—all of which are important for survival [140, 141]. The evolutionary hypothesis is further augmented by the observed co-development of abilities to detect and ignore irrelevant or unreliable vocal signals [142, 143]. A case in point is the common observation that urban humans today seem to have lost the acuity with which they note and differentiate natural sounds—such as those that may indicate wind speed, approaching weather conditions etc. Presumably this is because the modern urban human is well protected against the vagaries of nature, and these sounds are no longer as relevant to their survival as they were in the past centuries. On the other hand, an urban driver may be acutely aware of even faint sounds made by their vehicle on the road as it runs over obstacles, and may be able to judge the nature of the obstacle by the sounds. This shift in auditory awareness also affects the judgments made of other people within a crowded urban environment. Often, people lose the ability to "tell" much about the new people they talk to each day. We must however note that while these observations are commonly made, and support the evolutionary hypothesis, formal studies have not been performed to establish the connections between living styles and loss of voice-based judgment in urban humans.

Of the voice signal characteristics that have been found to be relevant to human judgment, some examples are tone and the sub-components of voice quality (explained in Chap. 6). Alterations in these have been observed to convey sarcasm, irony, disapproval, contempt and other elements of interaction [144–146]. Alterations in *prosody* have been related to various affective cues in vocal interactions [147]. *Speaking rate* has been related to the judgment of competence [148, 149] and other personality characteristics [150, 151]. *Loudness* has been reported to convey perceptions of dominance and dynamism, among other things [152]. *Pitch variations* are thought to convey impressions of benevolence [153]. These are just a few examples. There is a wealth of information in the literature about signal-level characteristics in relation to the judgments people make about other people. Later chapters in this book give more details.

### 1.3.2 Reactions Evoked by Voice

Voice evokes reactions in the speaker as well as in the listeners. Hearing oneself speak activates a feedback mechanism that is important for proper production and delivery of speech. It modifies the activity in the brain's auditory complex. Hampered audiovisual feedback mechanisms have been implicated in voice disorders such as aphasia, stuttering and schizophrenic voice hallucinations [154]. Often actors report that even the emotions that they merely portray through voice modify their mood, breathing, heart-rate etc. A very interesting perspective on the human voice is given

in [155], in which audible expression or vocalization of expression is analyzed in this context.

In listeners, voice evokes many kinds of reactions such as changes in feelings, mood and emotions. These in turn trigger physiological (such as hormone surges), psychological and behavioral changes. The first recorded study that attempted to categorize such effects was carried out in the 2nd century AD by Julius Pollux, a Greek grammarian, scholar and rhetorician from Naukratis, Egypt. According to historical accounts, Emperor Commodus appointed him a professor-chair of rhetoric in Athens at the Academy—*on account of his melodious voice*! Julius Pollux was the author of "Onomasticon," which still survives, and is, according to the Encyclopedia Britannica, "a Greek dictionary in ten books, each dedicated to Commodus, and arranged not alphabetically but according to subject-matter." His study of voice (vocal quality) from the listeners' perspective is reported to be part of these books.

The reactions evoked by voice in humans (and animals) continue to be studied and recorded. It is well known today, for instance, that the sound of music evokes activity in the brain, and consequently evokes emotions in listeners [156]. The same is true for the sound of human voice, and is supported by a multitude of studies e.g. [157–161]. Characteristics such as speech rate and loudness have been found to affect psychological states such as fear, sadness, anxiety, depression etc. [162]. Amplitude, frequency, pitch range, speed, inharmonicity etc. have been found to be used in conveying a sense of urgency [163]. A large number of reactions have been shown to be evoked by alterations in voice quality, prosody and energy. The *amygdala*, the brain's integrative center for emotions situated to the anterior of the hippocampus, is affected significantly by voice [164, 165], as is the pre-frontal cortex [166] and other areas of the brain, as these studies indicate. The amygdala also controls emotional behaviors, such as motivation. In fact, the net effect of voice on the brain is profound in many ways, and the pathways to evoking reactions in the limbic and nervous systems are well established today. In the light of these, it is only reasonable to believe that the features in voice that carry the necessary cues, if found, could indicate the presence of changes (such as hormone levels) that are brought about by the known reactions evoked by them. This chain of causes and effects is very relevant to profiling.

## 1.4   Computational Profiling

Computational profiling—what this book is largely about—is essentially an informed transfer of the process of making judgments about humans into the computing domain. From the studies mentioned above and in the literature, it is clear that human judgments are constrained by the limitations of human hearing, the brain's interpretative abilities, and the physical and mental states of the listener. As a result, they are often inaccurate. They are also inconsistent—they may be different at different times, although the voice signal continues to have the same information content.

Machines, of course, do not have these drawbacks. Once the process is computationalized, most ambiguities of repeatability and accuracy that play up in humans,