

Automatic Speaker Profiling from Short Duration Speech Data

Shareef Babu Kalluri^a, Deepu Vijayaseenan^a, Sriram Ganapathy^b

^a*Department of Electronics and Communication Engineering, National Institute of Technology
Karnataka, Surathkal, Mangalore, Karnataka, 575025, India*

^b*Learning and Extraction of Acoustic Patterns (LEAP) lab, Department of Electrical
Engineering, Indian Institute of Science, Bangalore, India.*

Abstract

Many paralinguistic applications of speech demand the extraction of information about the speaker characteristics from as little speech data as possible. In this work, we explore the estimation of multiple physical parameters of the speaker from the short duration of speech in a multilingual setting. We explore different feature streams for age and body build estimation derived from the speech spectrum at different resolutions, namely - short-term log-mel spectrogram, formant features and harmonic features of the speech. The statistics of these features over the speech recording are used to learn a support vector regression model for speaker age and body build estimation. The experiments performed on the TIMIT dataset show that each of the individual features is able to achieve results that outperform previously published results in height and age estimation. Furthermore, the estimation errors from these different feature streams are complementary, which allows the combination of estimates from these feature streams to further improve the results. The combined system from short audio snippets achieves a performance of 5.2 cm, and 4.8 cm in Mean Absolute Error (MAE) for male and female respectively for height estimation. Similarly in age estimation the MAE is of 5.2 years, and 5.6 years for male, and female speakers respectively. We also extend the same physical parameter estimation to other body build parameters like shoulder width, waist size and weight along with height on a dataset we collected for speaker profiling. The duration analysis of the proposed scheme shows that the state of the art results can be achieved using only around 1 – 2 seconds of speech data. To the best of our knowledge, this is the first attempt to use a common set of features for estimating the different physical traits of a speaker.

Keywords: Speaker Profiling, Short Duration, Formants, Harmonics.

1. Introduction

Apart from the textual message, human speech contains information about speaker identity, emotion, gender, accent etc. The extraction of speaker traits (parameters) from the speech data could further aid in speaker identification systems as well as in speaker clustering and diarization systems. The main challenge in estimating any such information is the separation of linguistic content and speaker traits.

In this paper, we try to address the problem of estimating physical parameters from the short duration of speech in a multilingual setting. This involves in predicting speaker meta information such as age and parameters of body build like height, weight, shoulder size and waist size. The motivation for height estimation range from biological understanding of the anatomy and its relationship to the speech properties to development of potential engineering systems for biometric applications [1, 2, 3]. While the current performance may not be applicable directly for developing robust solutions, the potential to augment speech based features as additional information has shown to improve other biometric methodologies based on finger printing [4]. In case of age estimation, researches have focused to identify the age group of a speaker (children, youth, adult and senior) from speech for most of the commercial applications (targeted advertisements, caller-agent pairing in call-centers etc) besides other applications like surveillance, forensics to narrow down on suspects from hoax/threat calls etc [1, 2, 5].

Speaker profiling is a challenging application area [6]. In many cases, there is no control over the amount of available speech data from the target speaker. Therefore, such systems are required to provide accurate predictions using a minimum amount of speech data. For example, DARPA RATS program targeted development of speaker and language recognition technology with as little as 3 seconds (s) of speech [7]. Thus, development of speaker profiling methods in short duration audio is important.

1.1. Physiological cues in speech

Literature shows that the physical dimensions of the speech production system are affected by the body build of a person. In general, a tall, well-built individual has lengthy vocal tract and large vocal folds [8]. The previous studies on the predicted height and weight of a person and their correlations with the acoustic features like fundamental frequency (F_0), vocal tract length (VTL) have generated

35 mixed results [9, 10, 11]. The correlation values of 0.53 (male) and 0.57 (female)
36 are reported between actual and perceived height values [10]. The previous studies
37 have also reported that VTL estimated from the speech has only a weak correlation
38 with body height [12, 13]. The only exception is a study [14] involving people
39 in the age group of 2.8 years to 25 years. This study reported the correlations
40 between actual vocal tract length and height using magnetic resonance imaging
41 (MRI). It shows that there is a strong correlation between vocal tract length and
42 height of the speaker for the subjects considered (0.88 for children, 0.83 for female
43 and 0.86 for male) [14]. It is also worthwhile noting that the sample size in this
44 study for adult subjects (17 to 25 years) was quite small (six female and 13 male).

45 One of the speech cues associated with the body size dimension of the speaker
46 is formant frequencies. They are weakly related to the body size dimensions such
47 as height and weight, and chest circumference [15, 16, 17]. The voice character-
48 istics of speech such as speech rate, sound pressure level, fundamental frequency,
49 etc. are affected by the speaker's age [18, 19, 20]. Other speech characteristics
50 like harmonics [21], jitter (micro variations in fundamental frequency), shimmer
51 (micro-variations of amplitude in frequency) occurs from age-related glottis
52 deterioration [22, 23] of the speaker. These features contain information about
53 speaker age.

54 Previous attempts [8, 10] in predicting the weight of a speaker, found a sig-
55 nificant correlation to exist between weight and vocal fold traits like dimensions
56 and mass. F_0 is significantly influenced by the obese and overweight people than
57 normal persons. The obese and overweight people have lower F_0 values than the
58 normal people [24]. A few studies show that the listeners are able to perceive the
59 weight (correlation of 0.724 for male and 0.627 for female speakers) and body
60 build [10, 25, 26]. Another study reports the correlation between log VTL and
61 log weight as 0.862, 0.875 and 0.903 for children, females and males respec-
62 tively [14]. While a weak correlation exists between the weight of the speaker
63 and the formant structure [15, 27], the speaking rate was found to be a useful
64 feature used by human listeners in weight attribute estimation [10].

65 While the past studies generate mixed results about the information present in
66 speech relating to speaker height, body dimensions and age, engineering appli-
67 cations to extract these physical traits from speech have shown practically useful
68 results (for example [28, 29]). However, in the existing literature, most of the
69 significant results have focused on the estimation of height and age from long
70 speech segments of few minutes ([29]) or by using hand labeled phoneme level
71 features [28]. The prior work on short duration speech shows that dealing with
72 utterances of 5sec. length is challenging yielding significantly worse results mak-

ing the systems inoperable for realistic applications [30]. In this work, we address the problem of reliably extracting height/age information from short duration speech 2 – 3sec. segments without using any phonetic information. We also extend the work to estimating more physical parameters (shoulder size, waist size, and weight). The main novelty of the proposed work lies in developing a unified framework for height/age and other physical parameter estimation. This is achieved using features that extract spectral structure of speech signal in terms of formant frequencies (peak locations in wide-band spectrum estimated using an autoregressive model) and harmonic frequency locations.

1.2. Organization of the Paper

The rest of the paper is organized as follows. Section 2 describes about the speaker profiling literature, motivation to carry out this work and contributions of the paper. Section 3 briefs about datasets, features extracted and regression technique used to estimate the physical parameters using speech data. Section 4 describes about the experiments conducted to estimate height and age of a speaker in monolingual setting using TIMIT dataset. Also this section discuss about the experiments performed on multiple physical parameters on multilingual setting using AFDS dataset in Section 4.3. A duration analysis is performed to know the minimum amount of speech data required for estimating physical parameters and this is explained in Section 4.4. Finally, conclusions are presented in Section 5.

2. Speaker Profiling Literature

While there is information about height/age in the speech signal, the extraction of these parameters is challenging, as these parameters are also affected by numerous other factors such as the content being spoken, emotion and mood of the speaker, gender of the speaker etc. These factors degrade the performance of the height and age estimation methods.

2.1. Height Estimation

The height of a speaker can be estimated by standard sound specific features such as formants, F_0 , sub-glottal resonances (SGR), short term spectral features and accumulated statistical features of the speech features across the sentence as a input to system.

The researchers predict the height of a speaker using the speech based features by using the short term features – Mel Frequency Cepstral Coefficients

(MFCC) [31, 32], Linear Prediction Coefficients(LPC) [31], formant frequencies [31, 33, 28], sub-glottal resonances [34, 35] and fundamental frequency [31]. Phone specific (vowels like /iy/, /ae/, /ey/, /ih/, /eh/ etc.) short term features like (MFCC, LPC) and formants shows a correlation of around 0.75 and for F_0 it is 0.59 in estimating the height [31]. In an alternate approach [36], the sub-glottal resonances are used for height estimation. SGRs are the resonance frequencies of sub-glottal (below the glottis) input impedance measurements from the top of the trachea. The SGRs are measured using the bark scale difference of the formants [35]. These resonances are shown to be correlated with the height information, and a simple polynomial relation can then be employed to estimate the height. Using the SGRs, the overall mean absolute error (MAE) of 5.4 cm, root mean square error (RMSE) of 6.8 cm at the sentence level and 5.3 cm, 6.6 cm of MAE and RMSE respectively at speaker level on TIMIT data.

A few other studies use the vowel regions (/aa/, /ae/, /ao/, /iy/) to predict the height of a person by formant track regression [28, 33]. This method obtained the MAE is reduced to 6.36cm for male and 6.8cm for female speakers by considering a subset of speakers and selected sentences from TIMIT dataset. By fusing the formant track regression with height distribution based classification, the MAE is 5.37cm and 5.49cm for male and female speakers respectively. Later line spectral frequencies were added to the feature set resulting in a lower MAE 4.93cm and 4.76cm for male and female speakers respectively. However, these approaches require speech transcription and phone level alignment.

Another set of approaches that do not depend on the speech transcriptions use accumulated statistics of the short term speech features as input. These features are typically used on a regression scheme (Support Vector Regression (SVR), Artificial Neural Networks (ANN), etc.) in predicting the height of a person. For example, various statistics like mean, median, percentiles etc. are extracted from the short-term spectral features for automatic height estimation [37, 38]. Here a set of features are selected from a large pool of statistical features. A feature selection algorithm precedes the support vector regression which provides the estimate of the height and obtains MAE of 5.3cm and RMSE of 6.8cm on TIMIT dataset. A similar approach uses i-vectors (dimension reduced version of background Gaussian Mixture Model (GMM) statistics) followed by regression schemes (SVR, ANN, etc.) to estimate the height of a speaker [3, 39].

In another approach, the height is divided into different bins and the height class of the speaker is estimated [32, 40]. For example the MFCC features are modeled using a background GMM to estimate the height class of a speaker (i.e., for a given utterance the height class is estimated). This approach using the TIMIT

dataset reports results with a RMSE of 6.4 cm and 5.7cm for male and female speakers respectively [40].

Singh et al. [41] reports that the MAE performance of the default predictor (average value of that parameter over the training set) is often better than the results in literature such as [33, 37, 38]. This study focuses on a bag of words representation instead of GMMs. The short term spectral features at multiple temporal resolutions are used to form a bag of words representation. For the TIMIT dataset, the MAE is 5.0 cm and RMSE is of 6.7 cm for male speakers and for female speakers the MAE is 5.0 cm and 6.1 cm RMSE. This study uses the short durations of speech data to estimate the height of a speaker [41].

2.2. Age Estimation

The accumulated statistics of the prosodic features and short term features can be used to estimate the age of the speaker. A popular approach uses prosodic features such as jitter / shimmer, harmonics to noise ratio, fundamental frequency [18, 22, 23]. These feature statistics are used by machine learning models like Artificial Neural Networks (ANN - Multilayer Perceptron), Support Vector Machines (SVM), k-Nearest Neighbor (KNN) etc. to classify the age group of a speaker. By considering both male and female genders the age class accuracy is 94.61% using an ANN model in proprietary dataset [18]. There have also been attempts to combine information from various levels such as short-term spectrum, prosodic features etc. These features are preceded by background GMM, SVM etc. for the age estimation [21, 23]. With Interspeech2010 Para linguistic challenge dataset, the unweighted accuracy was 52% and weighted accuracy was 49.5% for the age classification problem [21]. However, these efforts do not estimate the age, but only classify the input speaker as belonging to one of the age groups (e.g., kid, young adult, adult, etc.).

The statistical approaches adapted by researchers for age-group classification are Gaussian Mixture Model (GMM) Universal Background Model (UBM) [22, 42, 43], support vector machines [44, 45, 46], ANN [39]. These are followed by the statistical representation of short term features like MFCC, LPC, Perceptual Linear Prediction (PLP) coefficients, Temporal Patterns (TRAPS) [43] etc. In another approach, the age of a speaker is estimated by using a bag of words representation in place of background GMM from short-term cepstral features. In this work, short duration of speech data was considered and obtained MAE of 5.5 years & RMSE of 7.8 years for male, and for female speakers, MAE is 6.5 years & RMSE is 8.9 years on TIMIT dataset [41].

Using the UBM based approach, the short-term features are represented as supervectors/i-vectors and these are used as input features to a classifier [45, 29, 47]. Using NIST SRE08 and SRE10 data, the fusion of different short term features and i-vectors results in MAE of 4.7 years for male with correlation of 0.89, female MAE is 4.7 years with correlation of 0.91 [29]. A more recent approach using the deep neural networks on the short utterances of telephone speech using long short term memory (LSTM) recurrent neural networks(RNN) [48] MAE and correlation of male and female speakers are 8.72y, 0.37, and 7.95y, 0.54 respectively when 3s of speech is considered. An end to end deep neural network architecture using the x-vectors has also reported recently. Using only x-vectors on end to end system the MAE, correlations for 5s chunks of speech data are 5.78y, 0.74 for male, 4.23y, 0.87 for female respectively [30]. Table 1 shows the summary of the prior works methods and features for height and age estimation tasks.

2.3. Other Physical Characteristics

There are very few studies to estimate the other parameters like weight, shoulder size, chest circumference, shoulder to hip ratio, smoking habits, etc.,

The body size parameters like weight, neck etc. are predicted using F_0 and formants of all the vowels. The correlation between F_0 and first four formants with weight is 0.3 for male speakers [15]. Another study [16] shows the correlations of average fundamental frequency with shoulder circumference ($r = -0.29$), chest circumference ($r = -0.28$), shoulder-hip ratio ($r = -0.49$) and weight with formants is ($r = -0.43$).

Using the i-vector frame work weight is estimated and obtained the correlation of 0.56 for male and 0.41 for female speakers. The smoking habits are also predicted by using the i-vector framework with a log-likelihood ratio cost of 0.81 [39].

2.4. Limitations of Prior work

Majority of the speaker profiling works of the past concentrate on estimating only one physical parameter – either age or height. The best results in height estimation are obtained by using features that are phoneme specific [28, 31, 33]. This comes with the constraint on the system to have accurate transcription of the speech utterances with phone level alignment. The approaches involving SGR features [34, 35, 36, 40] require a separate dataset to learn the relationship between speech formants and the sub-glottal resonances. Other literature, often report the results on longer speech utterances using NIST recordings ($> 10s$) [3, 29, 30,

Table 1: Summary of prior work in age and height estimation.

Literature summary on Age			
Reference	Motivation	Features	Model
[18, 22, 23, 42]	Target advertisements	Pitch, jitter, shimmer, MFCC, LPC, etc.	ANN/SVM/GMM and fusion
[29, 45, 47]	Forensics, target advertisements	i-vectors	SVM / SVR
[30, 48]	Forensics, target advertisements, commercial applications	i-vectors/ x-vectors	DNN
[21, 43, 46]	Target advertisements	MFCC, Prosodic features, Formants, Pitch, PLPs, TRAPs	SVM/ GMM .
Literature summary on Height			
[3]	Forensics, biometric applications	i-vectors	LSSVR/ANN
[37, 38]	Forensics, biometric applications	OpenSmile	SVR
[28, 32, 33]	Forensic, biometric applications	LSF,Formants,MFCC	Linear Regression, GMM
[34, 35, 36, 40]	Relation between SGR and height	SGR	GMM, polynomial regression
Literature summary on Height and Age			
[39]	Forensics, target advertisements	i-vectors	LSSVR/ANN
[41]	Forensics, target advertisements	Short term spectral features	Random Forest

216 45, 47, 48] and does not address speaker profiling from short utterances. Even
217 for the i-vector based systems, the i-vectors may not be well estimated for short
218 utterances [29, 45, 47]. Also often gender specific speaker profiling results are
219 not reported [31, 38] and it was later reported that the gender-wise results of these

220 methods are inferior to default predictor based on the mean of the training data
221 performance genderwise [41]. So far the only work that addressed both height and
222 age estimation from short duration speech is Singh et.al. [41]. However, the prior
223 work on short duration speech shows that dealing with utterances of < 5 sec. of
224 speech in physical parameter estimation is challenging. To the best of the authors’
225 knowledge, literature does not address the physical parameter estimation from
226 short duration multilingual speech data.

227 2.5. *Contributions from This Work*

228 In this work, we attempt to address the main two challenges for physical trait
229 estimation, one is short duration of utterances and second is the multilingual na-
230 ture of the data. The goal is to come up with a common feature input for all phys-
231 ical parameter prediction systems. The proposed features do not require phone
232 level transcriptions. We consider different characteristics of the speech signal –
233 short-term spectral features, fundamental frequency, formant frequency locations
234 and narrow-band speech harmonics. With many experimental results, we show
235 that the proposed approach of using spectral features is useful in the prediction of
236 height/age and other physical attributes of the speaker.

237 To the best of our knowledge, this paper presents the first work of its kind
238 to illustrate the estimation of physical parameters from short durations of speech
239 signal in a multilingual setup. We perform height and age estimation experiments
240 in the TIMIT database [49] where the speech recordings are 2 – 3 seconds du-
241 ration. The combination of these features attain the MAE of 5.2years (male) and
242 5.6years (female) in age estimation and in case of height estimation the MAE is
243 of 5.2cm for males and 4.8cm for female speakers. In these experiments, the com-
244 bination of proposed features shows significant improvements over the previously
245 published results on the same dataset [35, 41]. We extend the same approach to
246 multilingual setting to predict multiple physical parameters like shoulder width,
247 waist size, weight along with height on a dataset. Finally, we investigate the mini-
248 mum amount of speech required to perform physical parameter estimation on both
249 TIMIT and AFDS datasets.

250 3. Methodology

251 In this work, we use two datasets for our experiments and analysis. One is
252 the standard TIMIT dataset [49], and the second one is a multilingual dataset,
253 Audio Forensic Dataset (AFDS) [50], collected for this purpose. We extract three
254 different features which doesn’t require the phoneme level transcriptions for short

Table 2: Statistics of each parameter in the TIMIT dataset [49]

Physical Characteristic	Minimum	Maximum	Mean	Standard Deviation
Male Speakers				
Height (<i>cm</i>)	157.48	203.20	179.73	7.09
Age (<i>y</i>)	20.63	75.77	30.52	7.57
Female Speakers				
Height (<i>cm</i>)	144.78	182.88	165.80	6.71
Age (<i>y</i>)	21.08	67.35	30.03	8.70
Male and Female Speakers				
Height (<i>cm</i>)	144.78	203.20	175.50	9.47
Age (<i>y</i>)	20.63	75.77	30.37	7.98

speech segments. The utterance level statistics of the extracted features is given to a support vector regression to estimate the physical parameters.

3.1. Datasets

The TIMIT dataset has 630 speakers, each speaker has contributed 10 recordings. Each of the ten recordings per speaker is considered as a separate input data sample. For training set, we have 462 speakers (326 male and 136 female speakers) and for testing 168 (56 female and 112 male speakers). The statistics of the dataset is given in Table 2. The training and validation splits has 4610 utterances which includes 3260 utterances from male speakers and 1360 utterances from female speakers. The test split has 1120 utterances from male speakers and 560 utterances from the female speakers. Each input utterance had 1 – 3 seconds of speech data for height and age prediction.

The second one is a dataset, collected from diverse dialects of individuals across India for this study. This dataset was named as Audio Forensic Dataset (AFDS) [50]. This dataset contains the speaker details like height, weight, shoulder width, waist size along with the speech utterances. Speech is recorded at a sampling frequency of 16 kHz. Each speaker provided around 2 minutes of speech data in three sessions, with each session lasting around 40 seconds. This

Table 3: Statistics of each parameter in the AFDS dataset [50]

Physical Characteristic	Minimum	Maximum	Mean	Standard Deviation
Male Speakers				
Height (<i>cm</i>)	156	188	171.0	6.7
Shoulder width (<i>cm</i>)	40	53	45.0	2.5
Waist size (<i>cm</i>)	68	112	86.0	7.6
Weight (<i>kg</i>)	45	107	67.9	11.1
Female Speakers				
Height (<i>cm</i>)	147	169	157.6	5.1
Shoulder width (<i>cm</i>)	30	45	38.4	2.6
Waist size (<i>cm</i>)	64	97	80.4	7.0
Weight (<i>kg</i>)	39	77	52.7	6.9
Male and Female Speakers				
Height (<i>cm</i>)	147	188	168.0	8.5
Shoulder width (<i>cm</i>)	30	53	43.5	3.7
Waist size (<i>cm</i>)	64	112	84.7	7.8
Weight (<i>kg</i>)	39	107	64.5	12.1

dataset is linguistically diverse with people having 12 different native languages. Each speaker is asked to read news articles in their native language as well as in English. This speech corpus contains 207 speakers including 161 males and 46 females. The speakers are in the age group of 18–37 years. The height, shoulder width and waist size are measured in centimeters (cm) and weight in kilograms (kg). The statistics of the dataset are tabulated in Table 3.

For evaluation purpose, the dataset is divided into training and testing datasets. Training data has 137 speakers (951 utterances) consisting of 104 males (727 utterances) and 33 females (224 utterances). Testing data has 70 speakers (538 utterances) consisting of 57 males (434 utterances) and 13 females (104 utterances). Train and test splits includes both English and native language. Both the training and testing splits contain speakers across the 12 different languages. There is no overlap of speakers in both the datasets of training and test splits.

286 3.2. Feature Extraction

287 In this paper, we try to come up with a common set of features that can be
 288 used for the physical parameters estimation. We explore different features which
 289 uncover the underlying the spectral structure of the speech signal to estimate the
 290 physical parameters. The short-term mel spectrogram captures the gross level
 291 spectral characteristics used in predicting height and age of a speaker [3, 5, 23, 31,
 292 39]. The fundamental and formant frequencies contain information about physical
 293 parameters of a speaker [15, 28, 35]. The narrowband spectral harmonics capture
 294 the fine spectral structure on a coarse temporal scale. The log harmonics are used
 295 in estimating the age and gender of a speaker [21]. We use both frequency and
 296 amplitude of the spectral peaks as harmonic features (to capture jitter and shimmer
 297 characteristics of speech).

298 3.2.1. Feature Extraction Using Mel Filter bank Features & UBM

299 *Mel Filter Cepstral Coefficients (MFCC):* The MFCC features are the most
 300 commonly representations used in speaker recognition. The MFCC features are
 301 have some information relating to the vocal tract length [22, 31]. In the past, the
 302 MFCC features and their statistics have been employed followed by the regression
 303 scheme for height and age estimation [21, 37, 38, 39]. In our work, we extract 20
 304 mel frequency cepstral coefficients (using a window length of 25 ms with a shift
 305 of 10 ms) along with delta and double delta features (yielding 60 MFCC features).

306 *Mel Filter bank features:* In our work, we use the logarithm of the mel spectral
 307 energy in short-term windows (25ms with a shift of 10ms) of the speech signal.
 308 The mel filter bank features are the short energy features computed prior to the
 309 Discrete Cosine Transform (DCT) in the MFCC feature computation. We extract
 310 40 mel filter bank features. The short spectral features contain the phonetic infor-
 311 mation as well as the speaker information. We adopt a supervector [51] approach
 312 which can summarize the gross spectral changes in order to normalize the effect
 313 of phonetic information in the short-term spectral representation.

314 315 *Statistical Representation:*

316 In order to form a background UBM model, a Gaussian Mixture Model is es-
 317 timated from short-term spectral features. Let \mathbf{x}_i and \mathbf{y}_i be input MFCC feature
 318 (i.e, $\mathbf{x}_i \in \mathcal{R}^{60}$) and mel-filter bank feature (i.e, $\mathbf{y}_i \in \mathcal{R}^{40}$) corresponding to frame i
 319 respectively. Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ represents the input MFCC feature vectors
 320 and $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T\}$ represent mel filter bank features for an input utter-
 321 ance with T frames. The diagonal covariance GMM -UBM is trained on MFCC

322 features. The GMM probability density is :

$$f_{UBM}(\mathbf{x}) = \sum_{j=1}^M w_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \mathbf{C}_j) \quad (1)$$

323 where \mathbf{x} , denotes input feature vector (MFCC) and $\boldsymbol{\mu}_j, \mathbf{C}_j$ represent the mean
 324 and the diagonal covariance matrix of the j^{th} GMM component with weight w_j
 325 respectively. The frame level first order statistics for a given frame i and each
 326 GMM component j is computed as:

$$\mathbf{f}_i^j = \mathbf{y}_i p(j|\mathbf{x}_i), \quad (2)$$

327 where the a-posteriori probabilities of a GMM component j is given by:

$$p(j|\mathbf{x}_i) = \frac{w_j \mathcal{N}(\mathbf{x}_i, \boldsymbol{\mu}_j, \mathbf{C}_j)}{\sum_{j=1}^M w_j \mathcal{N}(\mathbf{x}_i, \boldsymbol{\mu}_j, \mathbf{C}_j)}. \quad (3)$$

328 We then concatenate all \mathbf{f}_i^j for all GMM components to obtain a super vector
 329 $\mathbf{F}_i = [f_i^1, f_i^2, \dots, f_i^j, \dots, f_i^M]$ which represents the utterance. The first order
 330 statistics for a given utterance is:

$$\mathbf{F} = \frac{1}{T} \sum_{i=1}^T \mathbf{F}_i \quad (4)$$

331 Intuitively, if each GMM component j corresponds to a different sound class, the
 332 average of \mathbf{f}_i^j over the frames i would represent the short-term spectral average of
 333 frames that belong to that sound class. These features are used in support vector
 334 regression to estimate the physical parameter.

335 3.2.2. *Extraction of Fundamental and Formant Frequency Features*

336 We compute the fundamental frequency from a wideband analysis of speech
 337 signal (temporal window size of 20ms with a shift of 10ms). The estimation is
 338 performed with the PEFAC algorithm [52] which combines noise rejection and
 339 normalization while ensuring temporal continuity in the estimates using dynamic
 340 programming. For physical parameter estimation, we use the statistics (mean,
 341 standard deviation and percentiles) of the time varying fundamental frequency
 342 computed over the given speech recording. The formant frequencies are estimated
 343 by picking the peaks of an auto regressive (AR) model of the power spectrum. The

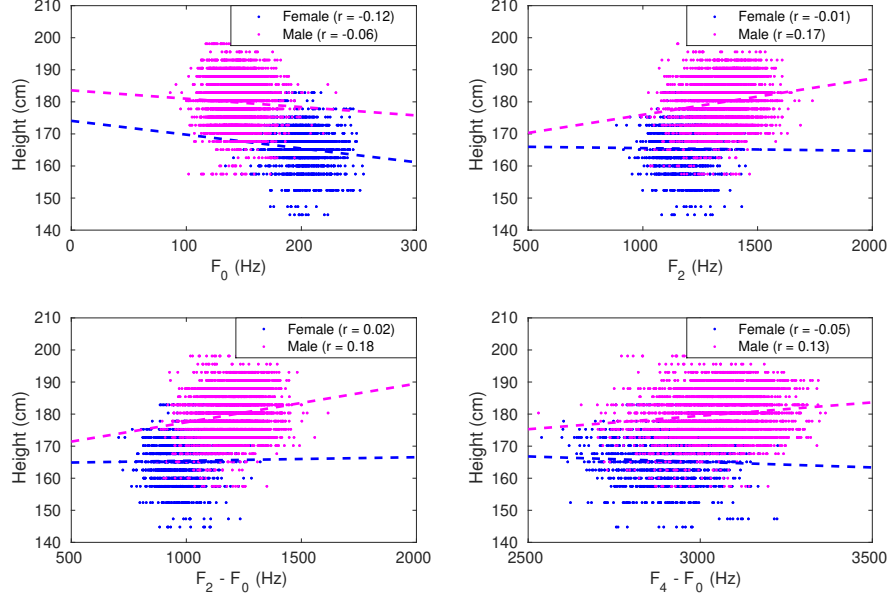


Figure 1: Scatter plot of fundamental and formant frequency estimates with the speaker height for TIMIT training set. Value in the brackets shows the correlation (r) between formants and corresponding physical parameter (height) for male and female speakers. The best fit line is also shown for both male and female speakers separately.

344 peaks of the wide-band (window length of $20ms$ with a shift of $10ms$) spectrum
 345 can approximately represent the formant structure. We use an AR model of order
 346 18 to extract peak locations results in nine peak locations. The first four peak
 347 locations are used to capture formant frequencies (denoted as F_1 , F_2 , F_3 and F_4).

348 The first four formant frequencies (F_1 , F_2 , F_3 , F_4) are extracted from the speech
 349 signal. We analyze the correlation between the fundamental frequency (F_0) and
 350 the other formant frequencies with the height values. The studies have shown F_0
 351 is inversely proportional to height of a speaker (indicating that the speakers with
 352 more height values have low fundamental frequency and vice-versa for speakers
 353 with lesser height values) [10, 16, 17]. The fundamental frequency (F_0), has a
 354 weak correlation with height ($r = -0.12$) for female speakers. Similarly, for
 355 male speakers F_2 showed a weak correlation with height value ($r = -0.17$).
 356 The correlations of male height vs F_0 ($r = -0.06$) and female height vs F_2
 357 ($r = -0.01$) are relatively modest. Literature has reported weak correlations
 358 between body build of the speaker and different functions of formant frequen-

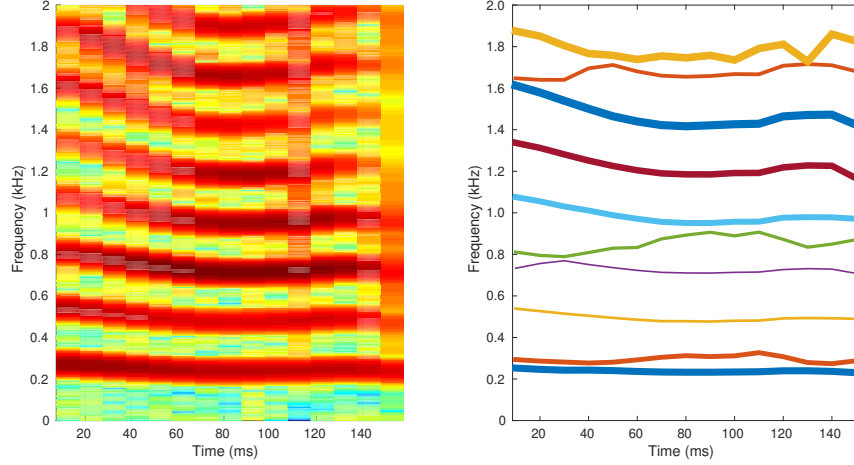


Figure 2: Spectrogram for vowel /AE/ and corresponding trajectories of first 10 peaks locations in a narrow-band spectrogram estimated using an AR model.

cies such as dispersion[53], average formant position[54], formant spacing [55], difference between F_0 and formants [15]. For example, we find the correlations between difference of F_0 and formants ($F_1 - F_0, F_2 - F_0, F_3 - F_0, F_4 - F_0$), Fig.1 depicts some of the results for the training portion of TIMIT dataset. It is observed that, $F_2 - F_0$ and $F_4 - F_0$ have weak positive correlation for male speakers ($r = 0.18$ and $r = 0.13$ respectively) and weak correlations for female speakers with height values [15].

Speaker identification systems have used mean value of pitch, range of pitch etc., as utterance level features [56]. In this work, we use a similar approach where each sentence is represented using statistics of the log fundamental frequency and log formant frequencies across the utterance. We use percentiles of log-peak locations in the short-term spectrum of speech (computed over time). The peak locations in the spectrum include the fundamental frequency and formant frequencies. In addition to the percentiles, the statistics of peak locations (in log-frequency scale) like the mean and standard deviation are used to estimate the physical parameters like height/age. These statistics can implicitly capture the average value, range and variance of fundamental frequency and formants.

3.2.3. Extraction of Harmonic Features

In addition to the conventional mel frequency spectrum and formants, we also experimented with the use of harmonic structure of the speech signal. The har-

monics are formed as a result of vocal fold vibration during voiced speech. It has been shown that variations in frequency (jitter) and amplitude (shimmer) contain useful information about age as well [22].

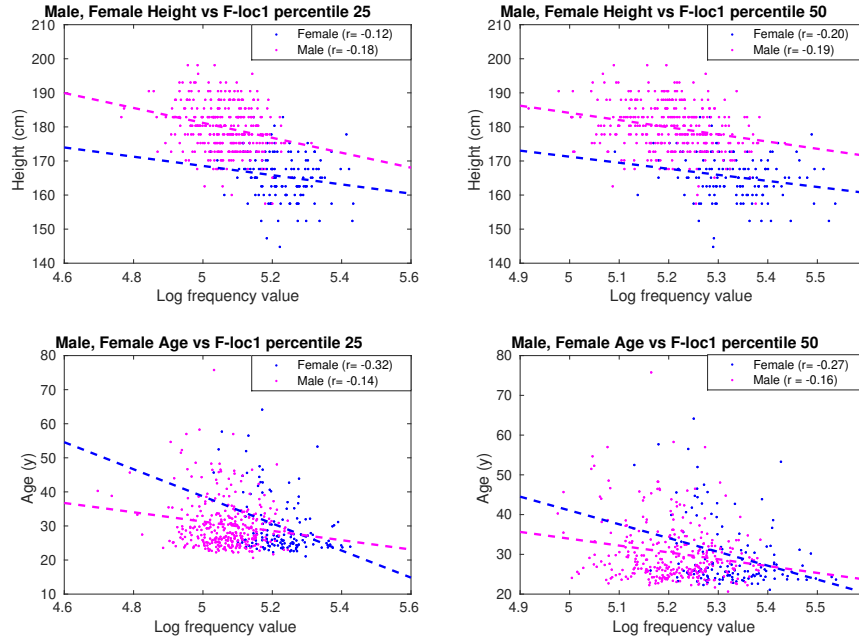


Figure 3: Scatter plot of Harmonic percentiles (25 and 50) vs physical parameter (height and age) for male and female speakers of TIMIT training data. Correlation (r) value between harmonic percentile and physical parameters (height and Age) is given in brackets for male and female speakers. The best fit line is also shown for both male and female speakers separately.

Using an AR model (order 80 with a window length of 60ms and a shift 10ms) of the spectrum, the peak locations (locations of the poles of the AR model) are identified. The logarithm of the frequency and amplitude of spectral peaks are computed at each frame. Each sentence is represented by the percentiles of log frequency and log amplitude values of spectral peaks over the utterance. The percentiles of harmonic frequencies represents the mean range and jitter in the harmonics. Similarly, the statistics on amplitude can contain shimmer in addition to average and range values. The collection of these statistics is referred to as “harmonic features” in this work. Fig.2 shows a short term spectrogram of the speech along with estimated harmonics.

The scatter plot for first harmonic frequency percentiles (25 and 50) on TIMIT

393 training data are shown in Fig.3 for both male and female speakers. It is observed
 394 that there is a weak negative correlation in case of height and age for percentiles
 395 25 and 50 for both male and female speakers. We also observe that the log mag-
 396 nitude statistics (percentiles) of the first two harmonic frequencies show a weak
 397 negative correlation with both age and height for both male and female speakers.
 398 These statistical harmonic features are used as input for support vector regression
 399 algorithm. The frequency location features capture jitter features and amplitude
 400 features captures shimmer features.

401 3.3. Prediction Using Support Vector Regression

402 We use a standard support vector regression (SVR) [57] as the model for pre-
 403 dicting the target of each physical parameter values given the input features. Let
 404 us denote the set of pair of input features along with target values as $\{(\mathbf{y}_1, t_1),$
 405 $(\mathbf{y}_2, t_2), \dots (\mathbf{y}_m, t_m)\}$. The function $f(\mathbf{y}) = \mathbf{w}^T \mathbf{y} + b$ corresponds to the linear
 406 SVR to learn and performs the following optimization:

$$\begin{aligned} \min \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ subject to} \\ & |\mathbf{w}^T \mathbf{y}_i + b - t_i| < \epsilon \end{aligned} \quad (5)$$

407 where b is the bias term and the “fit” function is controlled by the parameter ϵ . The
 408 maximum deviation from the target values is ϵ . The SVR optimization function
 409 aims to reduce the deviation from the target values by the parameter ϵ . We have
 410 also explored both linear and nonlinear kernels in this paper. In case of multiple
 411 features, we average the individual SVR outputs.

412 4. Experiments and Results

413 We perform height and age estimation experiments on TIMIT dataset. We use
 414 the standard train and test split in TIMIT. The algorithms are benchmarked using
 415 Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) metrics.

$$\begin{aligned} MAE &= \frac{1}{N} \sum_i |x_i^{pred} - x_i^{true}| \\ RMSE &= \sqrt{\frac{1}{N} \sum_i (x_i^{pred} - x_i^{true})^2} \end{aligned} \quad (6)$$

416 where x_i^{pred} and x_i^{true} are the predicted and target values for the i^{th} test utterance.

4.1. Results with Individual Features

In order to understand the effect of each feature separately, we evaluated the individual performance of the features. All hyper parameters of the system (e.g., kernel choice for SVR) and the order of the models were fixed based on the validation dataset performance.

We first perform a speech activity detection [58] and then extract the speech features. In order to extract the first order statistics (Fstats), we first train a 256 component GMM with 60 dimension MFCC features (\mathbf{x}_i). The Fstats are computed with 40 dimensional mel filter bank features (\mathbf{y}_i) using the Eq. 4. This gives $40 * 256 = 10240$ dimensional vector. The Fstats are fed to a support vector regression model to predict the physical parameters. A linear kernel is used for the support vector regression.

Fundamental frequency and formant features are extracted by picking the resonant frequencies of an all-pole model. A 18^{th} order (fixed based on validation set) model is used with a $20ms$ length window with $10ms$ shift. The 5^{th} , 25^{th} , 50^{th} , 75^{th} and 95^{th} percentile values across the entire utterance are employed as features. A linear kernel is used in the SVR.

A similar approach was followed in case of harmonic features. Thirty harmonics were extracted from an 80 order all-pole model, computed over a longer time window (length $60ms$ and shift $10ms$). The same set of percentiles are computed and used as input to a SVR with a third degree polynomial kernel (the order, window size and kernel are fixed based on the validation dataset). We separately evaluate the performance of harmonic frequencies, amplitudes as well as both together.

For comparison purposes, we also compute the Training data Mean Predictor (TMP). This just corresponds to providing the sample mean of the training data targets (physical parameters) as the estimate for any input, i.e., without using any evidence from the test speech. Fig.4 illustrates the performance of each feature as well as the TMP. In addition to the Fstats, and formants features, the figure also illustrates the effect of estimated harmonic frequency locations (F-loc) and corresponding amplitudes (Amp) as well as their combination ('harmonic' features). Both formants and Fstats have shown minimal improvement over TMP for both the genders in estimating the height of a speaker. The harmonic features show improvements only for female height and age estimation. In both these cases, the combination of harmonic features performs better than using either frequency locations or amplitudes. The performance improvement over TMP MAE is of 2.71% when Fstats are used for predicting height of male speakers. Similarly, for female speakers the improvement in MAE is of 4.01%, 3.23% , and 3.13% when

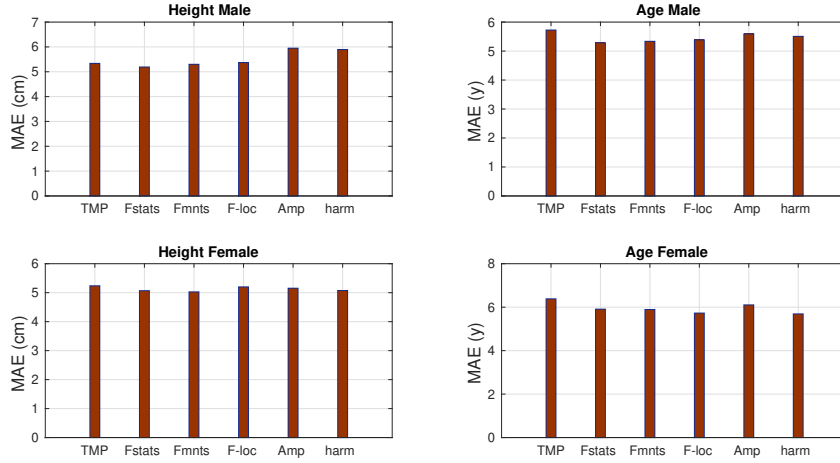


Figure 4: Mean absolute error comparison with training mean predictor (TMP) and prediction of different systems using first order statistics(Fstats), formants (Fmnts), harmonic frequency locations (F-loc), amplitude (Amp) and harmonic features (harmonic frequency locations & amplitude features together: harm) for height (left side) and age (right side) estimation using the TIMIT dataset.

formants, Fstats and harmonics are used respectively. For in predicting the age, all the features have shown a better performance when compared with TMP MAE for both the genders. For the male speakers, the improvement in MAE is of 6.8%, 3.82% and 7.7% for formants, harmonics and Fstats respectively. Similarly, for female speakers the improvement in MAE is of 7.71% 10.85% and 7.38% when formants, harmonics and Fstats respectively.

4.2. Results with Feature Combination

In our analysis, we found that the different feature sets produce different height and age estimation errors for a large number of validation speakers. With this knowledge, we attempt a simple averaging of the individual regression outputs to improve the final height and age estimates. We have made three different sets of feature combinations of Fstats and formant features with either harmonic frequency location (Comb -1) or amplitude (Comb -2) or harmonic features (both frequency and amplitude features: Comb -3). All our analyses use the standard TIMIT train and test splits. Table 4 reports the results along with the recent baseline which uses the standard train and test splits of TIMIT dataset [41].

The relative improvement of height prediction MAE for Comb-3 w.r.t TMP is 1.89% and 8.33% for male and female speakers respectively. Similarly, the

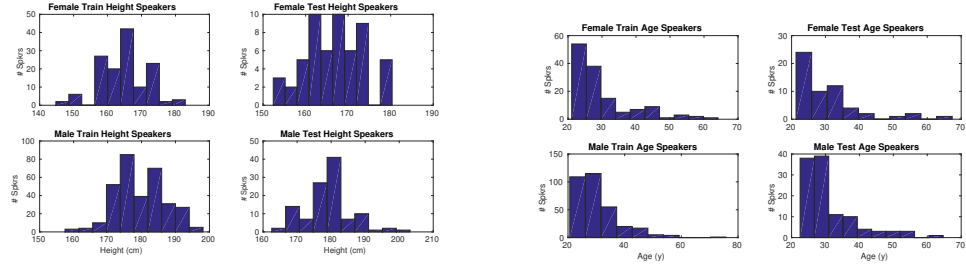
Table 4: Comparison of the proposed feature combinations – Comb -1 (Fstats + formant + frequency locations), Comb -2 (Fstats + formant + amplitude), Comb -3 (Fstats + formant + harmonic features (amplitude + frequency locations)) with state-of-the-art results on TIMIT dataset.

Height (cm) Estimation						
	Male		Female		All	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
TMP	5.3	7.0	5.2	6.5	7.4	9.0
[38]	-	-	-	-	5.3	6.8
[35]	5.6	6.9	5.0	6.4	5.4	6.8
[41]	5.0	6.7	5.0	6.1	-	-
Comb-1	5.2	6.8	5.0	6.3	5.2	6.8
Comb-2	5.2	6.9	4.8	6.2	5.2	6.7
Comb-3	5.2	6.8	4.8	6.1	5.2	6.7

Age(y) Estimation						
	Male		Female		All	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
TMP	5.7	8.1	6.4	9.2	5.9	8.4
[41]	5.5	7.8	6.5	8.9	-	-
Comb-1	5.3	8.2	5.8	9.2	5.5	8.7
Comb-2	5.3	8.2	5.6	8.8	5.4	8.6
Comb-3	5.2	8.1	5.6	8.7	5.4	8.5

relative improvement of age prediction MAE is 8.77%, and 14.29% for male and female speakers respectively. In case of RMSE, the relative improvement in height prediction of Comb-3 w.r.t to TMP is 2.94% and 6.15% for male and female speakers respectively. Similarly, for age prediction there is an 5.75% relative improvement for female speakers and no improvement for the male speakers.

We performed a paired t-test comparing the absolute errors from proposed system (Comb -3) and the default predictor (TMP) in a gender-wise manner. For both the tasks of height and age estimation, the proposed system is significantly different from the TMP ($p < 0.05$) across both the gender cases.



(a) Speaker Height – Training data (Left) and Test data (Right) (b) Speaker Age – Training data (Left) and Test data (Right)

Figure 5: Histogram of TIMIT dataset gender specific Training data and Test data – height and Age.

In case of height estimation, we also compare with three other baselines. The error metrics MAE and RMSE of the proposed systems as well as the baseline results are presented in Table 4. In case of female speakers both MAE and RMSE performances of Comb -3 are better than the baseline for height estimation. In order to gain further insight into the proposed height estimation system, we analyze the performance of height and age estimation of the data in different subgroups of Comb -3.

Table 5 lists various subgroups along with the height estimation performance and number of training speakers in each subgroup. It can be seen that large errors occur for speakers in the sub groups which are at the two extreme height values (row 3 & 6 for male speakers and 2 & 5 for female speakers) in Table 5. This may be due to the small amount of training data available for these groups. The gender specific histogram of speaker heights for both training and testing datasets are depicted in Fig.5a. We also observe that there is a mismatch in train and test height histograms. Such mismatches could have also resulted large error in extreme values of height.

Table 5: Height (h) estimation errors (MAE and RMSE in centimeters(cm))across different height subgroups using TIMIT test data

		Male			Female		
Sl. No.	Range	# Train Spkrs	MAE	RMSE	# Train Spkrs	MAE	RMSE
1.	$145 \leq h < 150$	0	-	-	2	-	-
2.	$150 \leq h < 160$	2	-	-	20	9.3	9.6
3.	$160 \leq h < 170$	15	11.9	12.2	75	2.5	3.0
4.	$170 \leq h < 180$	137	4.7	5.7	35	6.4	7.1
5.	$180 \leq h < 190$	140	2.9	3.7	3	14.9	14.9
6.	$190 \leq h < 203$	32	12.5	13.1	0	-	-

Table 6: Age (a) estimation error (MAE and RMSE in years) across different age subgroups using TIMIT test data

		Male			Female		
Sl. No.	Range	# Train Spkrs	MAE	RMSE	# Train Spkrs	MAE	RMSE
1.	$20 \leq a < 25$	67	4.6	4.8	47	2.7	3.0
2.	$25 \leq a < 30$	132	1.8	2.1	46	2.0	2.4
3.	$30 \leq a < 35$	66	2.9	3.4	14	4.7	5.2
4.	$35 \leq a < 40$	28	7.8	8.1	9	8.8	8.9
5.	$40 \leq a < 45$	13	13.0	13.1	9	13.0	13.1
6.	$45 \leq a < 55$	16	22.2	22.4	7	24.9	25.0
7.	$55 \leq a < 65$	3	35.5	35.5	3	21.9	21.9
8.	$65 \leq a < 76$	1	-	-	0	35.0	35.1

498 In case of age estimation, the only work that has reported results on short seg-
499 ments in TIMIT is by Singh et al. [41]. Comparison of this baseline with our
500 results and TMP is presented in Table 4. Note that in case of female speakers the

baseline had a higher MAE as compared to TMP. The proposed systems outperforms the baseline results and TMP in terms of MAE for male and female speakers. However, RMSE value is at par with TMP in case of Comb -3 male speakers and better than state of the art in female speakers in all the feature combinations. We analyzed the performance of Comb -3 for age estimation system by dividing the data into different subgroups as shown in Table 6. The RMSE is high over the TMP is due the presence of last three age groups (from 45 years to 75 years) in both the genders (refer Table 6). All these age groups have very few training speakers. Therefore, the RMSE error in these three groups are large (greater than 22y) and is dominates the overall RMSE performance. The histogram of gender specific speaker age in both training and testing datasets are depicted in Fig.5b. It can be seen that there are very few number of speakers above 45 years in training.

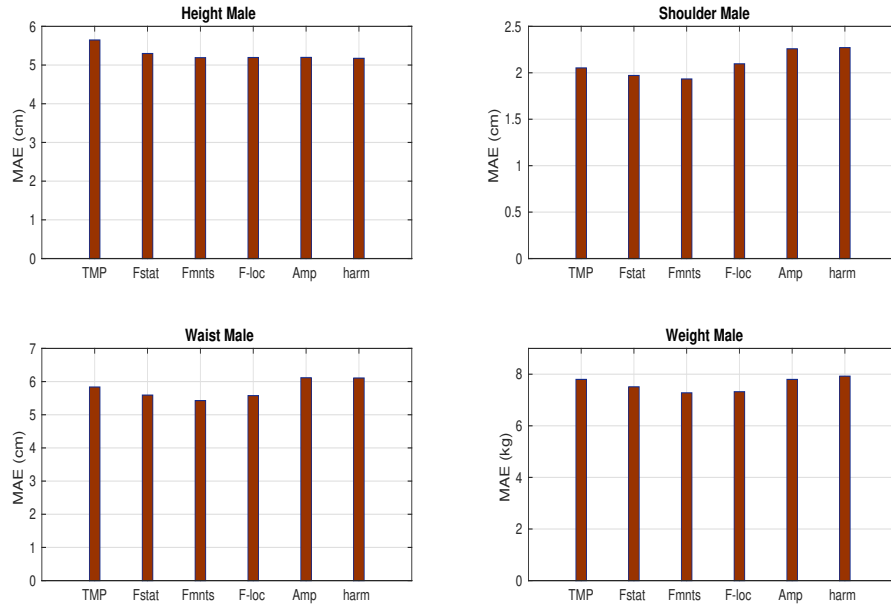


Figure 6: Mean absolute error of male speakers compared with training mean predictor (TMP) and prediction of different features i.e, first order statistics(Fstats), formants (Fmnts), harmonic frequency locations (F-loc), amplitude (Amp) and harmonic features (harmonic frequency locations & amplitude features together: harm) of physical parameters (Height, Shoulder width, Waist and Weight) using AFDS.

4.3. Extension to other physical parameters

We extend the same approach followed to estimate height and age to more physical parameters in a multilingual setting using the AFDS dataset as described in Section 3.1. We have analyzed the correlation of height with other parameters like shoulder size, waist size and weight on AFDS dataset. In the case of height the correlation values are small (0.2, 0.3 and 0.4 for shoulder size, waist size and weight respectively) for male speakers. The correlation values with age was negligible. Thus, these are parameters that cannot be predicted from height or age. We do not report results on only female data since, the number of female speakers is small.

In this regard, we use the same feature set (i.e, fundamental frequency, formants, harmonic features, and first order statistics of the speech signal) as explained in Section 3.2. In order to compute the first order statistics on AFDS, we have extracted 20 MFCCs along with deltas and double deltas and 40 filter bank features. We have used the GMM UBM learned from training data of TIMIT dataset itself, as the number of training utterances are less in AFDS. The Fstats are computed on AFDS using the Eq.4 (refer to Section 3.2.1).

The fundamental frequency, formants and harmonic features are extracted from the AFDS speech data along with its percentiles as explained in Section 3.2.2 and Section 3.2.3. These statistical features are fed to the support vector regression for the physical parameter estimation. The mean absolute error of each feature is compared with the training data mean predictor of each physical parameter (height, shoulder size, waist size and weight) is shown in Fig.6. The Fstats and formants shows better MAE performance for all the physical parameters. The harmonic features are better than TMP in case of height estimation.

Simple averaging is then performed on the predicted test targets obtained from formant features, Fstats and harmonics features (refer to Section 4.2) . The comparison of combination results and training data mean predictor are listed in Table 7. The table also lists an earlier algorithm developed by the authors as the baseline [50]. All the results use the same train and test split described in Section 3.1 (same splits are used in our previous work [50]). The baseline performs support vector regression of a bag of words representation extracted from the short-term spectrum of the speech. The performance metrics both MAE and RMSE on Comb -3 are better than the baseline for all speakers (both male and female speakers) except in MAE of weight estimation. With Comb -3, there is a substantial improvement of MAE and RMSE in all the physical parameters when compared with the TMP when only male speakers are considered. For further analysis we use Comb -3 set of features.

Table 7: Comparison of the proposed feature combinations – Comb -1 (Fstats + formant + frequency locations), Comb -2 (Fstats + formant + amplitude), Comb -3 (Fstats + formant + harmonic features (amplitude + frequency locations)) with baseline results of AFDS

Multiple Physical parameter Estimation – All (Male + Female)										
	TMP		Baseline[50]		Comb-1		Comb-2		Comb-3	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
Height(cm)	6.8	8.2	5.2	6.6	5.1	6.3	5.0	6.1	5.0	6.1
Shoulder(cm)	2.8	3.4	2.1	2.6	2.0	2.4	2.0	2.4	1.9	2.4
Waist(cm)	5.6	7.3	5.4	7.1	5.3	6.9	5.4	6.9	5.5	7.0
Weight(kg)	8.3	10.57	6.7	8.9	6.9	9.0	7.0	8.9	6.9	8.8

Multiple Physical parameter Estimation – Male									
	TMP		Comb -1		Comb -2		Comb-3		
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	
Height(cm)	6.4	6.9	5.1	6.3	5.1	6.2	5.0	6.1	
Shoulder(cm)	2.1	2.5	2.0	2.4	2.0	2.4	2.0	2.4	
Waist(cm)	5.8	7.3	5.4	7.0	5.6	7.1	5.5	7.1	
Weight(kg)	7.8	9.6	7.3	9.2	7.4	9.2	7.4	9.1	

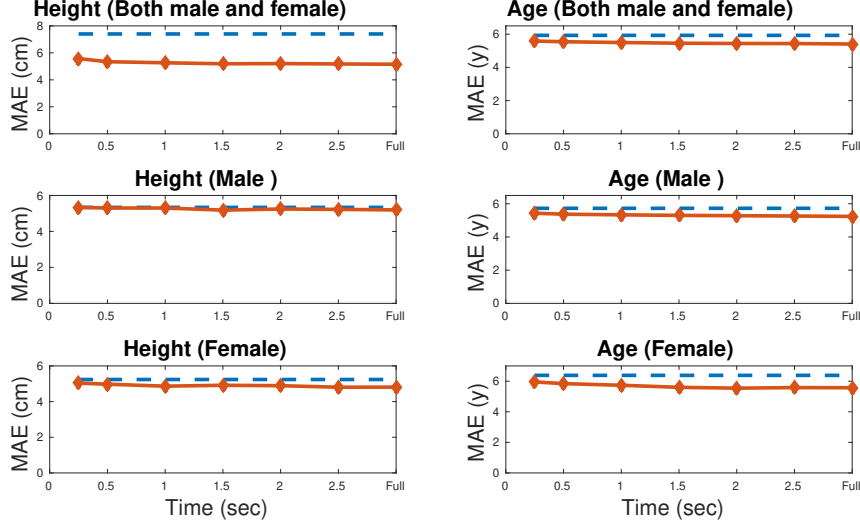


Figure 7: MAE vs duration of utterance, for physical parameters' (Height, Age) estimation from TIMIT database. The horizontal dashed line represent training data mean predictor (TMP) benchmark.

4.4. Duration Analysis

In order to analyze the minimum amount of speech required for the task, we try to evaluate the performance of the system at different utterance durations. We initially use the standard TIMIT database and evaluated the system for different time lengths of input speech ranging from 0.25s to full length. The mean absolute errors for these different lengths of speech were compared with TMP with height and age of a speaker and shown in Fig.7.

We performed a genderwise paired t-test comparing the absolute errors from proposed system (Comb -3) and the default predictor (TMP) for different durations of speech data. We find that (with criterion of $p < 0.05$) the proposed approach results in significant improvements in age estimation for all durations considered (starting from 0.5sec.) for both the genders and the relative improvement in MAE is 3.15% for males and 15.84% for female speakers. In the case of height estimation, the proposed approach results in significant improvements starting from 1.5 sec. duration of audio segments and the relative improvement in MAE for male speakers is 2.87% and for female speakers is 5.58%. Also, as the duration of the available speech increases, the MAE reduces as expected. Sub-

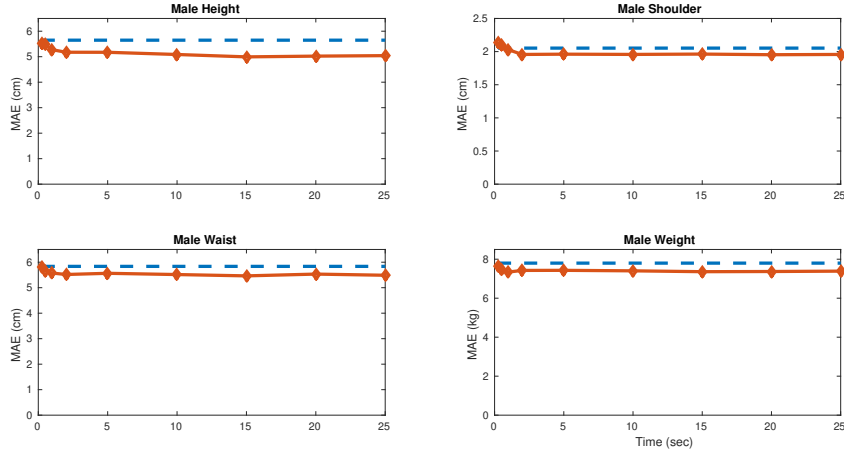


Figure 8: MAE vs duration of utterance, for physical parameters' (Height, Shoulder width, Waist size and Weight) estimation from AFDS database. The horizontal dashed line represent training data mean predictor (TMP) benchmark.

sequently, when sufficient amount of speech data is available, the mean absolute error get saturated.

It can be noted that even with roughly 1s of speech data, when both male and females speakers are considered, the model is able to obtain prediction error MAE of 5.27cm at par with Ganchev et al. [38] in speaker height prediction. As the available speech duration increases, this prediction error saturates around 5.2 cm when both genders are considered. Similarly for age prediction when both male and female speakers are considered together, the minimum duration of speech required to get the state-of-the-art prediction error (i.e, 5.5 years [38] is 0.5s. Even with around 3s speech available, the prediction error is marginally better (5.41 years). Gender wise results on duration analysis are also shown in Fig.7. About, 2s of speech data is required to get a performance comparable to the full length data.

We also extend the same duration analysis on other physical parameters (shoulder width, waist size and weight along with height) using the AFDS dataset. The system performance is evaluated for different lengths of speech files ranging from 0.25s to full duration(around 40s). We observed that mean absolute errors of each physical parameter for different durations' of speech signal is less than the training data mean prediction error except shoulder size by only using 0.5s for male speakers. From this, it is evident that the system is reliably able to predict the

physical parameters from 0.5s duration of speech signal with prediction error less than the training data mean. The duration of the speech at which the prediction error saturates is around 2s when both genders data is considered together. The mean absolute error for height is 5.1cm, shoulder width is 1.9cm, waist size is 5.4cm and for weight is 6.9 kg when we have 2s of speech data, where as when the available speech data is 40s, we have 5.0 cm, 1.9cm, 5.5cm and 6.9kg for height, shoulder width, waist size and weight respectively when both male and female speakers are considered together. The variation of MAE with respect to utterance duration for male speakers is shown in Fig.8. For male speakers also the MAE saturates around 2s as like above mentioned case (both male + female speakers). The change in MAE when full duration (40s) and 2s considered is 0.1cm in height, and there is no change in MAE for other physical parameters like shoulder size, waist size and weight estimation.

4.5. Summary

In short, it can be seen that each of the physical parameter prediction error is less than the TMP even with short speech segments (around 0.5 seconds). We are able to achieve the state-of-the-art results with around 1 – 2 seconds for all the physical parameters. The MAE of the proposed height estimation system on TIMIT (5.2cm for male, 4.8cm for female) is similar to the best height estimation results (5.0cm for male and 4.8cm for female) [28]. Note that this system [28] requires speech transcription for computing the phoneme specific features. In case of age estimation, the MAE of the proposed system (5.2 years for male and 5.6 years for female) is better than the state of the art result (MAE of 5.5 years for male, 6.5 years for female) reported on TIMIT [41]. Also, we demonstrate similar performances for other physical parameters in a multi-lingual setting. In summary, we hypothesize that the proposed methods could be used for speaker profiling where the duration of available speech data is limited.

5. Conclusions

In this work, we have explored the estimation of multiple physical parameters from short duration speech segments. In addition to conventional short-term spectral features, we also show that formant frequency features and harmonic structure of speech could be used as input to these tasks. Each of the individual features perform equally well on the test data and are able to achieve results that are comparable to state-of-the-art. Furthermore, these individual features are shown to be complementary and a simple averaging improves the performance by achieving an

MAE of 5.2 cm for male and all (male and female) and 4.8 cm for female speakers in height estimation. For age estimation, the MAE is 5.2 years, 5.6 years and 5.4 years for male, female and all speakers using the TIMIT dataset.

We have also presented the details of a new dataset where more speaker attributes like height, shoulder width, waist size and weight are collected. Each individual feature – first order statistics, formants, and harmonics – is able to achieve a prediction error less than the training data mean predictor in terms of MAE. The simple averaging of these predicted targets provides the best results in these tasks as well. While the proposed features and modeling are simple, we show that proposed approach is effective in various of speaker trait estimation tasks and outperform previously published results in these domains. To the best of authors knowledge, this is the first attempt to address the multilingual setting for speaker profiling tasks using short durations of speech data.

The duration analysis reveals that the prediction error of each physical parameter of a speaker is less than the training data mean predictor with as little speech as 0.5s. Also with around 1 – 2 seconds of data the MAE obtained is as good as the state-of-the-art results which were achieved using full duration of audio signal ($> 10s$). This enables the system to be useful in speaker profiling, speaker recognition tasks, targeted advertisements in commercial applications with short audio recordings from the target speaker. The extension to noisy speech in conversational setting would be the next logical step to developing forensic speech applications.

Acknowledgments

This work was partially funded by Science and Engineering Research Board (SERB) under grant no: EMR/2016/007934.

The authors would like to acknowledge the contribution of Sarthak Agrawal and Sanmathi Kamath who implemented the early version of the height/age prediction system while they were interning at the LEAP lab in Indian Institute of Science.

References

1. Nolan, F.. Forensic speaker identification and the phonetic description of voice quality. In: *A figure of speech: A festschrift for John Laver*. Psychology Press; 2005:385–411.

- 656 2. Singh, R., Keshet, J., Hovy, E.. Profiling hoax callers. In: *2016 IEEE*
657 *Symposium on Technologies for Homeland Security (HST)*. IEEE; 2016:1–6.
- 658 3. Poorjam, A.H., Bahari, M.H., Vasilakakis, V., et al. Height estimation from
659 speech signals using i-vectors and least-squares support vector regression.
660 In: *2015 38th International Conference on Telecommunications and Signal*
661 *Processing (TSP)*. IEEE; 2015:1–5.
- 662 4. Jain, A.K., Ross, A., Prabhakar, S., et al. An introduction to biometric
663 recognition. *IEEE Transactions on circuits and systems for video technology*
664 2004;14(1).
- 665 5. Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller,
666 C., Narayanan, S.. Paralinguistics in speech and language—state-of-the-art
667 and the challenge. *Computer Speech & Language* 2013;27(1):4–39.
- 668 6. Tanner, D.C., Tanner, M.E.. Forensic aspects of speech patterns: voice
669 prints, speaker profiling, lie and intoxication detection. Lawyers & Judges
670 Publishing Company; 2004.
- 671 7. Walker, K., Strassel, S.. The rats radio traffic collection system. In:
672 *Odyssey*. 2012:291–297.
- 673 8. Layer, J., Trudgill, P.. 1. phonetic and linguistic markers in speech. *Trudgill//Social Markers in Speech, Cambridge, CUP* 1979;(1-):1–32.
- 674 675 9. Gonzalez, J.. Estimation of speakers’ weight and height from speech: A
676 re-analysis of data from multiple studies by lass and colleagues. *Perceptual*
677 *and motor skills* 2003;96(1):297–304.
- 678 10. Van Dommelen, W.A., Moxness, B.H.. Acoustic parameters in speaker
679 height and weight identification: sex-specific behaviour. *Language and*
680 *speech* 1995;38(3):267–287.
- 681 11. Collins, S.A.. Men’s voices and women’s choices. *Animal behaviour*
682 2000;60(6):773–780.
- 683 12. Necioglu, B.F., Clements, M.A., Barnwell, T.P.. Unsupervised estimation
684 of the human vocal tract length over sentence level utterances. In: *2000*
685 *IEEE International Conference on Acoustics, Speech, and Signal Process-*
686 *ing. Proceedings (Cat. No. 00CH37100)*; vol. 3. IEEE; 2000:1319–1322.

- 687 13. Pisanski, K., Fraccaro, P.J., Tighe, C.C., O'Connor, J.J., Röder, S.,
688 Andrews, P.W., Fink, B., DeBruine, L.M., Jones, B.C., Feinberg, D.R..
689 Vocal indicators of body size in men and women: a meta-analysis. *Animal*
690 *Behaviour* 2014;95:89–99.
- 691 14. Fitch, W.T., Giedd, J.. Morphology and development of the human vo-
692 cal tract: A study using magnetic resonance imaging. *The Journal of the*
693 *Acoustical Society of America* 1999;106(3):1511–1522.
- 694 15. Rendall, D., Kollias, S., Ney, C., Lloyd, P.. Pitch (f 0) and formant profiles
695 of human vowels and vowel-like baboon grunts: The role of vocalizer body
696 size and voice-acoustic allometry. *The Journal of the Acoustical Society of*
697 *America* 2005;117(2):944–955.
- 698 16. Evans, S., Neave, N., Wakelin, D.. Relationships between vocal character-
699 istics and body size and shape in human males: an evolutionary explanation
700 for a deep male voice. *Biological psychology* 2006;72(2):160–163.
- 701 17. Greisbach, R.. Estimation of speaker height from formant frequencies. *In-*
702 *ternational Journal of Speech Language and the Law* 2007;6(2):265–277.
- 703 18. Müller, C.. Automatic recognition of speakers' age and gender on the basis
704 of empirical studies. In: *Ninth International Conference on Spoken Lan-*
705 *guage Processing*. 2006:.
- 706 19. Schötz, S.. Acoustic analysis of adult speaker age. In: *Speaker Classifica-*
707 *tion I*. Springer; 2007:88–107.
- 708 20. Schötz, S., Müller, C.. A study of acoustic correlates of speaker age. In:
709 *Speaker Classification II*. Springer; 2007:1–9.
- 710 21. Li, M., Han, K.J., Narayanan, S.. Automatic speaker age and gender
711 recognition using acoustic and prosodic level information fusion. *Computer*
712 *Speech & Language* 2013;27(1):151–167.
- 713 22. Müller, C., Burkhardt, F.. Combining short-term cepstral and long-term
714 pitch features for automatic recognition of speaker age. In: *Eighth Annual*
715 *Conference of the International Speech Communication Association*. 2007:.
- 716 23. van Heerden, C., Barnard, E., Davel, M., van der Walt, C., van Dyk, E.,
717 Feld, M., Müller, C.. Combining regression and classification methods for

- improving automatic speaker age recognition. In: *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE; 2010:5174–5177.
24. Souza, L.B.R.D., Santos, M.M.D.. Body mass index and acoustic voice parameters: is there a relationship? *Brazilian journal of otorhinolaryngology* 2018;84(4):410–415.
25. Lass, N.J., Brown, W.S.. Correlational study of speakers’ heights, weights, body surface areas, and speaking fundamental frequencies. *The Journal of the Acoustical Society of America* 1978;63(4):1218–1220.
26. Lass, N.J., Scherbick, K.A., Davies, S.L., Czarnecki, T.D.. Effect of vocal disguise on estimations of speakers’ heights and weights. *Perceptual and motor skills* 1982;54(2):643–649.
27. González, J.. Formant frequencies and body size of speaker: a weak relationship in adult humans. *Journal of phonetics* 2004;32(2):277–287.
28. Hansen, J.H., Williams, K., Bořil, H.. Speaker height estimation from speech: Fusing spectral regression and statistical acoustic models. *The Journal of the Acoustical Society of America* 2015;138(2):1052–1067.
29. Sadjadi, S.O., Ganapathy, S., Pelecanos, J.W.. Speaker age estimation on conversational telephone speech using senone posterior based i-vectors. In: *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE; 2016:5040–5044.
30. Ghahremani, P., Nidadavolu, P.S., Chen, N., Villalba, J., Povey, D., Khudanpur, S., Dehak, N.. End-to-end deep neural network age estimation. *Proc Interspeech 2018* 2018;;277–281.
31. Dusan, S.. Estimation of speaker’s height and vocal tract length from speech signal. In: *Ninth European Conference on Speech Communication and Technology*. 2005:.
32. Pellom, B.L., Hansen, J.H.. Voice analysis in adverse conditions: the centennial olympic park bombing 911 call. In: *Proceedings of 40th Midwest Symposium on Circuits and Systems. Dedicated to the Memory of Professor Mac Van Valkenburg*; vol. 2. IEEE; 1997:873–876.

- 749 33. Williams, K.A., Hansen, J.H.. Speaker height estimation combining gmm
750 and linear regression subsystems. In: *2013 IEEE International Conference*
751 *on Acoustics, Speech and Signal Processing*. IEEE; 2013:7552–7556.
- 752 34. Arsikere, H., Leung, G.K., Lulich, S.M., Alwan, A.. Automatic height
753 estimation using the second subglottal resonance. In: *2012 IEEE Interna-*
754 *tional Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
755 IEEE; 2012:3989–3992.
- 756 35. Arsikere, H., Leung, G.K., Lulich, S.M., Alwan, A.. Automatic estimation
757 of the first three subglottal resonances from adults’ speech signals with appli-
758 cation to speaker height estimation. *Speech Communication* 2013;55(1):51–
759 70.
- 760 36. Arsikere, H., Lulich, S.M., Alwan, A.. Automatic estimation of the
761 first subglottal resonance. *The Journal of the Acoustical Society of Amer-*
762 *ica* 2011;129(5):EL197–EL203.
- 763 37. Mporas, I., Ganchev, T.. Estimation of unknown speaker’s height from
764 speech. *International Journal of Speech Technology* 2009;12(4):149–160.
- 765 38. Ganchev, T., Mporas, I., Fakotakis, N.. Audio features selection for auto-
766 matic height estimation from speech. In: *Hellenic Conference on Artificial*
767 *Intelligence*. Springer; 2010:81–90.
- 768 39. Poorjam, A.H., Bahari, M.H., et al. Multitask speaker profiling for esti-
769 mating age, height, weight and smoking habits from spontaneous telephone
770 speech signals. In: *Computer and Knowledge Engineering (ICCKE), 2014*
771 *4th International eConference on*. IEEE; 2014:7–12.
- 772 40. Arsikere, H., Lulich, S.M., Alwan, A.. Estimating speaker height and sub-
773 glottal resonances using mfccs and gmms. *IEEE Signal Processing Letters*
774 2013;21(2):159–162.
- 775 41. Singh, R., Raj, B., Baker, J.. Short-term analysis for estimating physical
776 parameters of speakers. In: *2016 4th International Conference on Biometrics*
777 *and Forensics (IWBF)*. IEEE; 2016:1–6.
- 778 42. Metze, F., Ajmera, J., Englert, R., Bub, U., Burkhardt, F., Stegmann,
779 J., Muller, C., Huber, R., Andrassy, B., Bauer, J.G., et al. Comparison
780 of four approaches to age and gender recognition for telephone applications.

- 781 In: *2007 IEEE International Conference on Acoustics, Speech and Signal*
782 *Processing-ICASSP'07*; vol. 4. IEEE; 2007:IV–1089.
- 783 43. Bocklet, T., Stemmer, G., Zeissler, V., Nöth, E.. Age and gender recog-
784 nition based on multiple systems-early vs. late fusion. In: *Eleventh Annual*
785 *Conference of the International Speech Communication Association*. 2010:.
- 786 44. Spiegl, W., Stemmer, G., Lasarczyk, E., Kolhatkar, V., Cassidy, A., Potard,
787 B., Shum, S., Song, Y.C., Xu, P., Beyerlein, P., et al. Analyzing fea-
788 tures for automatic age estimation on cross-sectional data. In: *Tenth Annual*
789 *Conference of the International Speech Communication Association*. 2009:.
- 790 45. Bahari, M.H., McLaren, M., Van hamme, H., Leeuwen, D.v.. Age esti-
791 mation from telephone speech using i-vectors. In: *Thirteenth Annual Con-*
792 *ference of the International Speech Communication Association*. 2012:.
- 793 46. Li, M., Jung, C.S., Han, K.J.. Combining five acoustic level modeling
794 methods for automatic speaker age and gender recognition. In: *Eleventh*
795 *Annual Conference of the International Speech Communication Association*.
796 2010:.
- 797 47. Shivakumar, P.G., Li, M., Dhandhanian, V., Narayanan, S.S.. Simpli-
798 fied and supervised i-vector modeling for speaker age regression. In: *2014*
799 *IEEE International Conference on Acoustics, Speech and Signal Processing*
800 *(ICASSP)*. IEEE; 2014:4833–4837.
- 801 48. Zazo, R., Nidadavolu, P.S., Chen, N., Gonzalez-Rodriguez, J., Dehak,
802 N.. Age estimation in short speech utterances based on lstm recurrent neural
803 networks. *IEEE Access* 2018;6:22524–22530.
- 804 49. Garofolo, J.S., Lamel, L.F., Fisher, W.M., Fiscus, J.G., Pallett, D.S..
805 Darpa timit acoustic-phonetic continous speech corpus cd-rom. nist speech
806 disc 1-1.1. *NASA STI/Recon technical report n* 1993;93.
- 807 50. Kalluri, S.B., Vijayakumar, A., Vijayasenan, D., Singh, R.. Estim-
808 ating multiple physical parameters from speech data. In: *Machine Learning*
809 *for Signal Processing (MLSP), 2016 IEEE 26th International Workshop on*.
810 IEEE; 2016:1–5.

- 811 51. Reynolds, D.. An overview of automatic speaker recognition. In: *Pro-*
812 *ceedings of the International Conference on Acoustics, Speech and Signal*
813 *Processing (ICASSP)(S. 4072-4075). 2002:.*
- 814 52. Gonzalez, S., Brookes, M.. Pefac-a pitch estimation algorithm robust to
815 high levels of noise. *IEEE/ACM Transactions on Audio, Speech, and Lan-*
816 *guage Processing* 2014;22(2):518–530.
- 817 53. Fitch, W.T.. Vocal tract length and formant frequency dispersion correlate
818 with body size in rhesus macaques. *The Journal of the Acoustical Society of*
819 *America* 1997;102(2):1213–1222.
- 820 54. Puts, D.A., Apicella, C.L., Cárdenas, R.A.. Masculine voices signal men’s
821 threat potential in forager and industrial societies. *Proceedings of the Royal*
822 *Society B: Biological Sciences* 2012;279(1728):601–609.
- 823 55. Reby, D., McComb, K.. Anatomical constraints generate honesty: acoustic
824 cues to age and weight in the roars of red deer stags. *Animal behaviour*
825 2003;65(3):519–530.
- 826 56. Peskin, B., Navratil, J., Abramson, J., Jones, D., Klusacek, D., Reynolds,
827 D.A., Xiang, B.. Using prosodic and conversational features for high-
828 performance speaker recognition: Report from jhu ws’02. In: *Acoustics,*
829 *Speech, and Signal Processing, 2003. Proceedings.(ICASSP’03). 2003 IEEE*
830 *International Conference on; vol. 4. IEEE; 2003:IV–792.*
- 831 57. Smola, A., Vapnik, V.. Support vector regression machines. *Advances in*
832 *neural information processing systems* 1997;9:155–161.
- 833 58. Tan, Z.H., Lindberg, B.. Low-complexity variable frame rate analysis for
834 speech recognition and voice activity detection. *IEEE Journal of Selected*
835 *Topics in Signal Processing* 2010;4(5):798–807.