

THE ROLE OF VOICE SOURCE MEASURES ON AUTOMATIC GENDER CLASSIFICATION

Yen-Liang Shue and Markus Iseli

University of California Los Angeles
Dept. of Electrical Engineering
405 Hilgard Ave., Los Angeles, CA 90095
yshue@ee.ucla.edu, iseli@ee.ucla.edu

ABSTRACT

Physiological properties of the glottis and the vocal tract change with age and gender. Since these changes are reflected in the speech signal, acoustic measures related to those properties can be helpful for automatic age and gender classification. In this paper, the focus is on automatic gender classification, which is implemented using Support Vector Machines (SVMs), using acoustic measures that are related to the voice source. Acoustic measures of the vocal tract and the voice source were extracted from 3880 utterances spoken by 205 male and 160 female talkers (aged 8 to 39 years old). Formant frequencies and formant bandwidths were used as vocal tract measures, and open quotient and source spectral tilt correlates were used as voice source measures. Results show that the addition of these measures can help to improve automatic gender classification results for most age groups.

Index Terms— Voice source, Age classification, Gender classification

1. INTRODUCTION

Gender-based differences in human speech are due in part to physiological differences such as vocal fold thickness or vocal tract length, and differences in speaking style. Physiological properties of the glottis and the vocal tract change with age and gender. Since these changes are reflected in the speech signal, acoustic measures related to those properties can be helpful for age and gender classification. Assuming the linear source-filter model of speech production [1], the contribution of acoustic measures to such classification can then be attributed to the voice source or vocal tract. To our knowledge, with the exception of fundamental frequency (F_0), there has been no study that has examined the role of measures related to the voice source on age and/or gender classification.

It is well known that F_0 values for male talkers drop during adolescence due to a lengthening and thickening of the vocal folds. F_0 for adult males is typically around 120 Hz, while F_0 for adult females is around 200 Hz [2]. This effect is mostly due to a lengthening and thickening of the male vocal folds.

It is also well known that, due to vocal tract length differences, adult males exhibit lower formant frequencies than adult females [2]. Interestingly, for preadolescent children, studies also found lower formant frequencies for boys compared to girls of ages 5-6 [3], 7-8 years [4], and ages 5, 7, 9, and 11 years (for Australian English) [5]. These findings imply that, overall, boys have larger vocal tracts than girls. In [6], statistical analysis of children speech confirmed that formant frequencies (F_1 , F_2 , F_3), and not F_0 , differentiate gender for children as young as four years of age, while formant frequencies

plus F_0 differentiate gender after 12 years of age. These findings lead to the conclusion that for preadolescent children, vocal tract measures play a bigger role for gender classification than the voice source measure F_0 . For adult speech, automatic gender classification has been presented in [7] which used linear predictive coding (LPC)-derived measures that represent the vocal tract.

In [8], changes in magnitude and variability of, among other measures, F_0 , formant frequencies, and spectral envelope are presented as a function of age for talkers from 5 to 50 years old. For F_0 , the study showed a drop between ages 12 and 15 for males and a drop of F_0 variation for all talkers between ages 5 and 15. Formant frequencies (F_1 , F_2 , F_3) decreased between ages 10 and 15, where formant frequencies of male talkers decreased faster and reached much lower absolute values than those of female talkers. The study showed that children younger than age 10 displayed greater spectral variability than adults.

In [9], we analyzed age, sex, and vowel dependencies, for talkers between the ages of 8 and 39, of the following three voice source measures: F_0 ; $H_1^* - H_2^*$, the difference of the first two source spectral harmonic magnitudes (related to the open quotient [10]);¹ and $H_1^* - A_3^*$, the difference of the first source spectral harmonic magnitude and the magnitude of the source spectrum at the frequency location of the third formant (related to source spectral tilt [10]). For male talkers, the results showed a drop of about 5 dB in $H_1^* - H_2^*$ around age 15 and a continuous decrease of $H_1^* - A_3^*$ between ages 8 and 39 by about 10 dB. For female talkers, the value of $H_1^* - H_2^*$ remained relatively unchanged between ages 8 and 39, whereas for $H_1^* - A_3^*$ a slight decrease by about 4 dB was shown. These developmental changes resulted in higher values of F_0 , $H_1^* - H_2^*$, and $H_1^* - A_3^*$ for adult female talkers compared to adult male talkers [12].

In this paper, acoustic measures from both the voice source and the vocal tract are used for automatic gender classification of 8 to 39 year olds. The vocal tract measures consist of formant frequencies and formant bandwidths, and the voice source measures used are F_0 , $H_1^* - H_2^*$, and $H_1^* - A_3^*$. Training and testing is done using support vector machines (SVM's). It is analyzed if voice source measures can improve automatic gender classification. Finally, the SVM classification results are compared with human perceptual classification tests.

2. SPEECH DATA

Speech recordings from five age groups, ages 8–9, 10–11, 12–13, 14–15 and 16–39 were taken from the CID database [13]. Each

¹The asterisk indicates a correction for the influence of vocal tract resonances [11].

recording was of the form “I say uh, bVt again”, where the target vowel ‘V’ was /ih/, /eh/, /ae/ or /uw/. The vowel /iy/ in ‘bead’ was also used. These utterances were spoken at the habitual speaking level and most talkers repeated the phrases twice. For the analysis, only the manually segmented target vowels were used. The distribution of talkers (males/females) and number of utterances per age group is listed in Table 1. The total number of male/female talkers is 205/160 and the total number of utterances is 3880.

Table 1. Distribution of gender and utterances for each age group.

Age group	males/females	No. of utterances
8-9	48/36	810
10-11	48/33	807
12-13	38/34	708
14-15	22/21	413
16-39	49/36	1142

3. METHODS

The acoustic measures used for gender classification were the first three formant frequencies (F_1 , F_2 , and F_3), the first two formant bandwidths (B_1 and B_2), and the measures related to the voice source F_0 , $H_1^* - H_2^*$, and $H_1^* - A_3^*$. The third formant bandwidth, B_3 , was not used due to its large variance. The formant frequencies and bandwidth values were estimated using the “Snack Sound Toolkit” software [14] with these settings: analysis window length of 25 ms, window shift of 1 ms and pre-emphasis factor of 0.96. F_0 was extracted using the STRAIGHT algorithm [15]. The magnitudes H_1 , H_2 , and A_3 were estimated from the spectrum using the values of F_0 and F_3 . Corrections [11], denoted by the asterisks, were made to these measures to remove the effects of the vocal tract. For each of the voice source measures, a first order Legendre polynomial was fitted to the raw values to obtain a measure of the mean and the slope (denoted by Δ) across the duration of the vowel.

Classification was done using a Support Vector Machine (SVM) classifier with the Radial Basis Function kernel. In this study, the LIBSVM toolkit [16] was used to train and test on vectors containing different combinations of acoustic measures extracted from the five analyzed vowels. For each classification experiment, 70% of the utterances, selected randomly, were used for training; the remaining utterances were used for testing. Five experiments were performed for each combination of acoustic measures and the average accuracy recorded.

For our perception test, four male subjects between the ages 26 and 39 volunteered. They were each presented with 100 utterances of the target words and had to decide between male or female voice. The target words were manually segmented out of the whole carrier phrase and were played back in random order with state of the art headphones. The distribution of male and female utterances per age group were: age group 8-9, 7 males/7 females (m/f); ages 10-11, 8/8; ages 12-13, 8/8; ages 14-15, 12/10; and ages 16-39, 15/17. The same perception tests were also performed using just the segmented vowel part of the target word.

4. RESULTS AND DISCUSSION

For this section, formant information (F_1 , F_2 , F_3 , B_1 , and B_2) will be denoted by FB .

4.1. F_0 and formants

As a first step, we analyzed the contribution to gender classification accuracy of only F_0 , only FB , and F_0 plus FB (labeled by $M0$). These measures are the most widely used in gender and age classification.

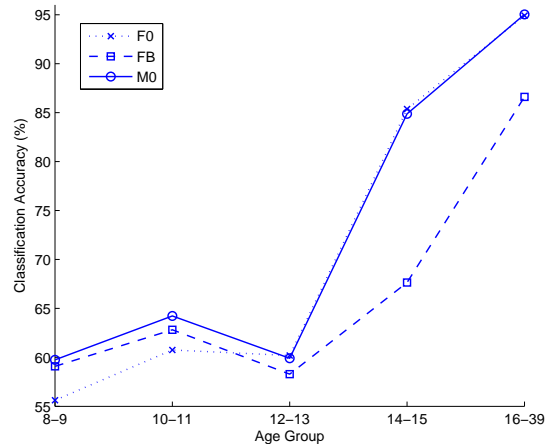


Fig. 1. Gender classification accuracy for each age groups using just F_0 , just FB , and F_0 plus FB ($M0$).

Figure 1 shows the classification accuracy for each age group using those measures. For ages 8 to 11 it can be seen that formant information only (FB) performs better than F_0 . This is consistent with [6]. Gender classification accuracy for ages 8 to 13 is always below 65%, but between age groups 12–13 and 14–15, it increases to 85% for F_0 and to 68% for FB ; these results can be attributed to the large drop of F_0 for males around ages 12 to 15 (about 105 dB on average) [9, 8] and to a decrease of formant frequencies for males relative to females [8]. Since $M0$ yielded the best results for all age groups, it was chosen as a baseline to compare the performance of voice source measures against. Table 2 shows the different measure sets.

Table 2. Measure sets ($M0$ - $M3$) used in the gender classification tests. $M0$, in bold, is used as the baseline measure set.

Set	Acoustic Measures					
	F_0	FB	$H_1^* - H_2^*$	$H_1^* - A_3^*$	ΔF_0	$\Delta H_1^* - H_2^*$
$M0$	✓	✓				
$M1$	✓	✓	✓			
$M2$	✓	✓	✓	✓		
$M3$	✓	✓	✓	✓	✓	✓

4.2. Adding voice source measures

Figure 2 compares the contribution in gender classification accuracy using the voice source measures to using the baseline measure set ($M0$), shown as a solid line. Table 3 shows the values corresponding to this figure. It can be seen that adding voice source measures plays a significant role only for age groups 10–11 and 12–13, where the absolute accuracy was improved by up to 9% using measure set $M3$. For age group 8–9 the accuracies are below 60% and the SVM seems

unable to separate the classes for males and females satisfyingly. Although it was shown in [9] that the source measures $H_1^* - H_2^*$ and $H_1^* - A_3^*$ are dependent on age and gender, the changes in classification accuracy for age groups 14–15 and 16–39 when using M_1 or M_2 are not significant. This could be attributed to the already large classification accuracy of the baseline.

A closer look at the classification accuracy results for age group 12–13 is shown in Table 4, which shows the percentage correct classification of males and females. Compared to M_0 , the addition of the voice source measures assists in increasing the classification accuracy by about 7% for males and 9% for females when using M_3 . However, since the measures of measure set M_2 are easier to calculate than those of M_3 , and M_2 showed a classification accuracy improvement for all ages between 10 and 39, it is recommended to use M_2 for gender classification. M_2 will be used throughout the remainder of this paper.

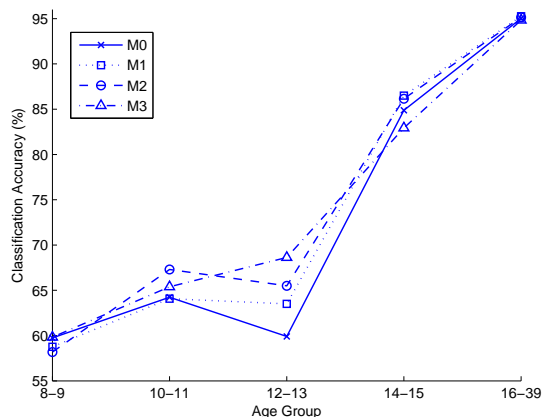


Fig. 2. Gender classification accuracy for each age group using the measures sets M_1 , M_2 and M_3 . M_0 represents the baseline performance results. These values are listed in Table 3.

Table 3. Gender classification accuracy results for the different measurement sets (M_0 - M_3) and age groups.

Age group	M_0 (Baseline)	M_1	M_2	M_3
8-9	59.75%	58.76%	58.18%	59.83%
10-11	64.23%	64.07%	67.30%	65.39%
12-13	59.91%	63.51%	65.50%	68.63%
14-15	84.88%	86.50%	86.18%	82.93%
16-39	95.03%	95.26%	95.15%	94.85%

4.3. Comparison with perception results

Table 5 compares automatic classification results (denoted by SVM) with human perception results from this study (denoted by OUR) and from perception experiments in Perry et al. [6] (denoted by PER). These perception experiments were done using the full target words. All values are gender recognition accuracies in percent. Dashes in the table represent unavailable data. SVM results were using measure set M_2 . The SVM classifier performs comparably with the human subjects for the talkers aged 14 and above. For talkers

Table 4. Gender classification results for age group 12-13, distinguishing between males and females classified correctly. Column for total

Set	M	F	Total
M_0	59.28%	60.60%	59.91%
M_1	63.24%	63.80%	63.51%
M_2	63.06%	68.20%	65.50%
M_3	66.67%	70.00%	68.63%

aged below 14, the results are somewhat mixed and the accuracies reduce with decreasing age; however this trend also exists with the human classifiers. In effect, in the “total” section of the table, the SVM results agree well with the perception results.

Since the SVM was only given short target vowels, and the listeners were able to listen to the whole target word, it seemed only fair to see how listeners would perform when given only short vowel segments. Interestingly, for talkers of age 15 and above, the results were similar to whole word gender classification (about 90% recognition accuracy). Our experimental subjects were mostly using F_0 to do the classification. For talkers of age 14 and below however, our experimental subjects agreed that their decisions were based purely on “guessing”.

The results show that using the voice source measures $H_1^* - H_2^*$ and $H_1^* - A_3^*$ with F_0 and formant information (i.e. measure set M_2) can help improve automatic gender classification for all age groups, except where the classifier was not able to sufficiently model the two gender classes; this occurred for the very young children (ages 8–9). This is consistent with the age and gender dependencies of these measures shown in [9].

Table 5. SVM gender classification results using measure set M_2 compared with perception results from this paper (OUR) and from Perry et al. [6].(PER). All values in percent. Dashes indicate unavailable values.

Age	8	9	10	11	12	13	14	15	16
Males									
SVM	-	-	-	-	67	-	83	-	94
OUR	39	-	-	72	91	-	100	-	100
PER	74	-	-	-	82	-	-	-	99.7
Females									
SVM	-	-	-	-	68	-	90	-	97
OUR	68	-	-	75	31	-	70	-	97
PER	56	-	-	-	56	-	-	-	95
Total									
SVM	58	-	-	67	66	-	87	-	95
OUR	54	-	-	73	61	-	86	-	98
PER	65	-	-	-	69	-	-	-	97

5. SUMMARY AND CONCLUSIONS

In this paper, we examined the role of voice source measures in automatic gender recognition and compared the results to perceptual experiments performed on the same database. Vocal tract and voice source measures were extracted from a large database of 3880 utterances spoken by 205 males and 160 females. For the vocal tract, the formant frequencies and formant bandwidths were used, while

for the voice source F_0 , $H_1^* - H_2^*$ (related to open quotient) and $H_1^* - A_3^*$ (related to spectral tilt) were used. The slopes (derivatives) were also calculated for the voice source measures. Automatic gender classification using SVMs was performed on five age groups with different sets of measures.

Using a baseline measure set consisting of F_0 , the first three formants and the first two bandwidths, it was found that the most consistent set of voice source measures was $H_1^* - H_2^*$ with $H_1^* - A_3^*$. For ages greater than 9, using these two measures increased the classification accuracy, although the improvements decreased for older talkers as the role of F_0 became more dominant. The measure sets which included the slopes ΔF_0 and $\Delta H_1^* - H_2^*$ did not produce consistent results and in some age groups actually reduced the classification accuracy.

With perception experiments which used whole words instead of only vowel segments, we were able to show that, with the exception of very young talkers, the results from the automatic classifier are quite comparable. This suggests that there exists cues outside the range of the vowels which could aid in automatic gender classification. Future work will focus on finding reliable methods to extract these cues.

6. REFERENCES

- [1] G. Fant, *Acoustic theory of speech production*, Mouton, The Hague, Paris, 1960.
- [2] G. E. Peterson and H. L. Barney, "Control methods used in a study of the vowels," *The Journal of the Acoustical Society of America*, vol. 24, no. 2, pp. 175–184, March 1952.
- [3] B. Weinberg and S. Bennett, "Speaker sex recognition of 5- and 6-year-old children's voices," *The Journal of the Acoustical Society of America*, vol. 50, pp. 1210–1213, 1971.
- [4] S. Bennett, "Vowel formant frequency characteristics of preadolescent males and females," *The Journal of the Acoustical Society of America*, vol. 69, pp. 231–238, 1981.
- [5] P. Busby and G. Plant, "Formant frequency values of vowels produced by preadolescent boys and girls," *The Journal of the Acoustical Society of America*, vol. 97, pp. 2603–2606, 1995.
- [6] T. L. Perry, R. N. Ohde, and D. H. Ashmead, "The acoustic bases for gender identification from children's voices," *The Journal of the Acoustical Society of America*, vol. 109, no. 6, pp. 2988–2998, June 2001.
- [7] K. Wu and D. G. Childers, "Gender recognition from speech. part i: Coarse analysis," *The Journal of the Acoustical Society of America*, vol. 90, no. 4, pp. 1828–1840, 1991.
- [8] S. Lee, A. Potamianos, and S. Narayanan, "Acoustics of children's speech: Developmental changes of temporal and spectral parameters," *The Journal of the Acoustical Society of America*, vol. 105, no. 3, pp. 1455–1468, March 1999.
- [9] M. Iseli, Y.-L. Shue, and A. Alwan, "Age, sex, and vowel dependencies of acoustical measures related to the voice source," *The Journal of the Acoustical Society of America*, vol. 121, no. 4, pp. 2283–2295, April 2007.
- [10] E. B. Holmberg, R. E. Hillman, J. S. Perkell, P. Guiod, and S. L. Goldman, "Comparisons among aerodynamic, electroglottographic, and acoustic spectral measures of female voice," *J. Speech Hear. Res.*, vol. 38, pp. 1212–1223, 1995.
- [11] M. Iseli and A. Alwan, "An improved correction formula for the estimation of harmonic magnitudes and its application to open quotient estimation," in *Proceedings of ICASSP*, Montreal, Canada, May 2004, vol. 1, pp. 669–672.
- [12] H. M. Hanson and E. S. Chuang, "Glottal characteristics of male speakers: Acoustic correlates and comparison with female data," *The Journal of the Acoustical Society of America*, vol. 106, pp. 1064–1077, 1999.
- [13] J.D. Miller, S. Lee, R.M. Uchanski, A.F. Heidebreder, B.B. Richman, and J. Tadlock, "Creation of two children's speech databases," in *Proceedings of ICASSP*, May 1996, vol. 2, pp. 849–852.
- [14] Kåre Sjölander, "Snack sound toolkit," KTH Stockholm, Sweden, 2004, <http://www.speech.kth.se/snack/> (last viewed Jan. 2007).
- [15] H. Kawahara, A. de Cheveign, and R. D. Patterson, "An instantaneous-frequency-based pitch extraction method for high quality speech transformation: revised TEMPO in the STRAIGHT-suite," in *Proceedings ICSLP'98*, Sydney, Australia, December 1998.
- [16] Chih-Chung Chang and Chih-Jen Lin, *LIBSVM: a library for support vector machines*, 2001, Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.