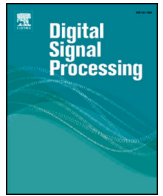




Contents lists available at ScienceDirect

Digital Signal Processing

www.elsevier.com/locate/dsp



Prominence features: Effective emotional features for speech emotion recognition

Shaoling Jing, Xia Mao, Lijiang Chen*

School of Electronic and Information Engineering, Beihang University, Beijing 100191, China

ARTICLE INFO

Article history:
Available online xxx

Keywords:
Prominence features
Speech annotation
Consistency assessment
Speech emotion recognition

ABSTRACT

Emotion-related feature extraction is a challenging task in speech emotion recognition. Due to the lack of discriminative acoustic features, classical approaches based on traditional acoustic features could not provide satisfactory performances. This research proposes a novel type of feature related to prominence, which, together with traditional acoustic features, are used to classify seven typical different emotional states. To this end, the author group produces a Chinese Dual-mode Emotional Speech Database (CDESD), which contains additional prominence and paralinguistic annotation information. Then, a consistency assessment algorithm is presented to validate the reliability of the annotation information of this database. The results show that the annotation consistency on prominence reaches more than 60% on average. Subsequently, this research analyzes the correlation of the prominence features with emotional states using a curve fitting method. Prominence is found to be closely related to emotion states, to retain emotional information at the word level to the greatest possible extent and to play an important role in emotional expression. Finally, the proposed prominence features are validated on CDESD through speaker-dependent and speaker-independent experiments with four commonly used classifiers. The results show that the average recognition rate achieved using the combined features is improved by 6% in speaker-dependent experiments and by 6.2% in speaker-independent experiments compared with that achieved using only acoustic features.

© 2017 Elsevier Inc. All rights reserved.

1. Introduction

In recent years, technology for automatic emotion recognition from speech has matured sufficiently to be applied in several real-life scenarios [1–3], such as call centers [4,5], auxiliary disease diagnosis [6–9], remote education [10–12], driving safety [13, 14], and computer games [15]. However, existing speech emotion recognition (SER) technology has not achieved very good performance, probably because of the lack of effective emotion-related features. To solve this problem, we explore the effect of prominence features in SER.

There is a diversity of viewpoints regarding the definition of prominence. Phonetician Zhao Li holds that prominence, simply speaking, refers to the fact that some syllables sound more prominent than the surrounding syllables in a speech stream, which is caused by different intensities that the speaker put on different syllables [16]. Streefkerk et al. regard prominence as the perceptual salience of a language unit. Sentence accent and pitch accent lead to perceived prominence, and there is no definite boundary

between sentence accent and pitch accent [17]. Meanwhile, Terken defines prominence as “words or syllables that are perceived as standing out from their environment” [18]. Kakouros et al. also agree with this definition [19]. In Chinese, a word is defined as a syllable. Traditionally, Chinese prominence has lexical stress and sentence accent. Hence, we use the term prominence here to refer to the accentuation of syllable within words or of syllables within sentences.

Speech prominence, which refers to a prosodic property, is typically viewed within the contexts of phonetic properties and physical realization [20]. In some languages, prominence can be a very subjective topic. If we observe the prominence trajectories of only a few speakers, it is likely that these prominence trajectories will vary greatly from one another. However, statistical methods are needed to study the subjectivity of prominence. In particular, the general idea is that the subjective perception of prominence is a response elicited to unpredictable prosodic trajectories in a normal train of speech, which draw the attention of the listener [19]. Indeed, Kakouros et al. have proven that the statistical distribution of prosodic cues can impact the subjective perception of word prominence [19]. Therefore, the subjectivity of prominence makes this research a challenging work. However, prominence has been

* Corresponding author.

E-mail address: chenlijiang@buaa.edu.cn (L. Chen).

<https://doi.org/10.1016/j.dsp.2017.10.016>

1051-2004/© 2017 Elsevier Inc. All rights reserved.

1 proven to be useful in many spoken language applications since it
2 encompasses affective expression and affects speech understand-
3 ing, especially in speech synthesis [21,22]. Researchers have proven
4 that integrating prominence patterns into speech recognition can
5 improve speech classifier performance [23,24]. However, so far, few
6 studies have applied prominence attributes to improve the accu-
7 racy of emotion recognition. Inspired by the use of prominence
8 by speakers to express their emotions or attitudes, we explore
9 whether prominence has an impact on SER. For this purpose, we
10 need to extract the prominence features from utterances. However,
11 the majority of studies on automatic prominence detection have
12 focused on supervised methods, which require manually annotated
13 prominence labels [25]. Therefore, the objective of our work is to
14 create prominence annotations, extract prominence features and
15 apply them in speech emotion recognition.

17 1.1. Related work

19 Generally, valid speech corpora are usually accompanied by
20 emotional labels, which are used to train models for SER. Most
21 of existing Chinese emotional speech databases have three main
22 problems. First, they are not available/public, thus not accessible
23 by researchers. Second, they are annotated for specific tasks, thus
24 do not contain enough information for general speech emotion
25 recognition task. Third, their annotations are confined to emo-
26 tional labels and transcription information, while corpora of other
27 languages contain much more rich annotation information. For
28 instance, the German database emoDB is annotated with pho-
29 netic transcriptions, stress, and segment and pause boundaries
30 [26]. The Japanese database UU is annotated with orthographic
31 transcriptions, dimensional emotions (such as pleasant–unpleasant
32 and aroused–sleepy) and paralinguistic information [27]. How-
33 ever, so far, the available Chinese emotional speech databases
34 lack such rich annotation information. For example, Chinese emo-
35 tional speech database DFEIC is only annotated with prosody, func-
36 tion and emotional states [28]. Paralinguistic information refers
37 to meaningful information, such as emotion or attitude, delivered
38 along with linguistic messages. Therefore, to study whether the
39 combination of paralinguistic information, prominence and sen-
40 tence function is useful for SER, the author group produces a new
41 Chinese Dual Modal Emotional Speech Database (CDESD), which
42 contains rich annotation information including text transcription,
43 syllable segment, initials, finals, tones, silence, voiced and unvoiced
44 segments, prominence, sentence function and emotional states. In
45 our emotional annotation, an utterance may express more than
46 one emotion, in that case the emotional state label will contain
47 all emotion states, unlike in the other databases.

48 Since follow-up experiments are based on the assumption that
49 the annotation information is correct, we must ensure that the la-
50 beled data is reliable before using them in SER. In light of the
51 appraisal theory from the domain of emotional psychology [29],
52 each annotator may have his or her own subjective perception of
53 the affective state expressed by an individual, thus the annotation
54 may be influenced by the identity of the annotator, annotator's
55 experience, memories, and reasoning, etc. [30,31]. Besides, the an-
56 notation task is very time-consuming and laborious; hence, the
57 long time required for annotation may affect the subjective per-
58 ceptions and judgments of the annotators. To ensure the reliability
59 of annotated data, we need to use relevant algorithms to deter-
60 mine the agreement among the annotators. Cohen's kappa [32] is
61 a measure of the agreement between two raters when determin-
62 ing to which categories a finite number of subjects belong. The
63 two raters either agree in their rating (i.e., the category to which a
64 subject is assigned) or they disagree. However, Cohen's kappa does
65 not account for the degree of disagreement. The weighted Cohen's
66 kappa measure [33] addresses this issue by using a predefined ta-

67 ble of weights that measure the degree of disagreement between
68 the two raters. However, neither method is suitable for our anno-
69 tation scheme. For the situation at hand, one of rater's jobs is to
70 determine the boundaries between different syllables. Since each
71 boundary's position can be selected from a continuous region, it
72 does not fall into a fixed number of classes. Besides, in some an-
73 notation layer, such as syllable and tones layer and initials, finals
74 and tones layer (see Section 3.1), each syllable may be divided into
75 uncertain number of classes and the division position is also uncer-
76 tain. Therefore, it is necessary to present a consistency assessment
77 algorithm specific to our situation.

78 After confirming the reliability of the prominence annotation,
79 we explore whether prominence can be used as a basis for ex-
80 tracting useful features to improve SER performance. Most previ-
81 ous studies have relied on acoustic features, such as pitch, en-
82 ergy and spectral coefficients [34–38]. However, in recent years,
83 many researchers have begun to incorporate additional informa-
84 tion to improve the accuracy of emotion recognition. Lee et al.
85 combined acoustic and language information for emotion recog-
86 nition, thereby improving negative emotion classification by 45.7%
87 and 32.9% compared with the use of only acoustic information
88 and only language information, respectively [39]. Schuller et al.
89 proposed to perform SER by combining acoustic features with lin-
90 guistic information [13]. Chen et al. proposed two classes of speech
91 emotional features extracted from electroglottography and speech
92 signals [40], and combined those features with traditional features
93 for SER [41]. Mencattini et al. proposed novel features related to
94 the amplitude modulation in a speech emotion prediction model
95 [42]. Moreover, paralinguistic and non-linguistic information [27,
96 43,44], emotional point information [45], and lexical and semantic
97 information [46] are also closely related to the emotional states
98 of human beings and play a crucial role in SER. In particular,
99 prominence is an important property of speech. Speakers may use
100 prominence to draw the listener's attention to specific points of an
101 utterance to express their emotions or attitudes. Previous works
102 have demonstrated the usefulness of prominence in clarifying am-
103 biguities in spoken utterances and constructing understandable
104 synthetic speech [47]. Therefore, it is feasible to use prominence
105 in SER system.

107 1.2. Main contributions

109 This research aims at producing a new speech emotion recog-
110 nition system based on traditional acoustic features and our
111 newly proposed prominence features. To this end, we identify
112 three aspects of the problem related to prominence features in
113 SER: (1) prominence annotation, (2) annotation consistency, and
114 (3) prominence effectiveness.

115 Fig. 1 presents the scheme of the proposed recognition sys-
116 tem. First, each syllable is annotated with its prominence, and a
117 consistency assessment algorithm is designed to validate the reli-
118 ability of the annotation. Second, the manually labeled prominence
119 information is extracted, the extracted prominence information
120 is pre-processed to obtain the prominence features, and a curve
121 fitting-based (CF-based) method is used to analyze the correla-
122 tion between the proposed features and emotional states. Third,
123 290 acoustic features are extracted, and the most useful features
124 are selected. The Sequence Backward Selection (SBS) algorithm
125 and the Sequence Forward Selection (SFS) algorithm are used for
126 feature selection. Fourth, the prominence features are combined
127 with acoustic features for use in classifying emotional states. The
128 Double Layer Feature Dimension Reduction (DLFDR) model, which
129 combines Principal Component Analysis (PCA) and Nonparametric
130 Discriminant Analysis (NDA) [48], is applied to reduce the fea-
131 ture dimensions. Finally, the novel features are tested in speaker-
132 independent and speaker-dependent experiments. Good perfor-

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات