

Multimodal sentiment analysis using hierarchical fusion with context modeling

N. Majumder^a, D. Hazarika^b, A. Gelbukh^a, E. Cambria^{*,c}, S. Poria^c

^a Centro de Investigación en Computación, Instituto Politécnico Nacional, Mexico

^b School of Computing, National University of Singapore, Singapore

^c School of Computer Science and Engineering, Nanyang Technological University, Singapore

ARTICLE INFO

Keywords:

Multimodal fusion
Sentiment analysis

ABSTRACT

Multimodal sentiment analysis is a very actively growing field of research. A promising area of opportunity in this field is to improve the multimodal fusion mechanism. We present a novel feature fusion strategy that proceeds in a hierarchical fashion, first fusing the modalities two in two and only then fusing all three modalities. On multimodal sentiment analysis of individual utterances, our strategy outperforms conventional concatenation of features by 1%, which amounts to 5% reduction in error rate. On utterance-level multimodal sentiment analysis of multi-utterance video clips, for which current state-of-the-art techniques incorporate contextual information from other utterances of the same clip, our hierarchical fusion gives up to 2.4% (almost 10% error rate reduction) over currently used concatenation. The implementation of our method is publicly available in the form of open-source code.

1. Introduction

On numerous social media platforms, such as YouTube, Facebook, or Instagram, people share their opinions on all kinds of topics in the form of posts, images, and video clips. With the proliferation of smartphones and tablets, which has greatly boosted content sharing, people increasingly share their opinions on newly released products or on other topics in form of video reviews or comments. This is an excellent opportunity for large companies to capitalize on, by extracting user sentiment, suggestions, and complaints on their products from these video reviews. This information also opens new horizons to improving our quality of life by making informed decisions on the choice of products we buy, services we use, places we visit, or movies we watch basing on the experience and opinions of other users.

Videos convey information through three channels: audio, video, and text (in the form of speech). Mining opinions from this plethora of multimodal data calls for a solid multimodal sentiment analysis technology. One of the major problems faced in multimodal sentiment analysis is the fusion of features pertaining to different modalities. For this, the majority of the recent works in multimodal sentiment analysis have simply concatenated the feature vectors of different modalities. However, this does not take into account that different modalities may carry conflicting information. We hypothesize that the fusion method

we present in this paper deals with this issue better, and present experimental evidence showing improvement over simple concatenation of feature vectors. Also, following the state of the art [1], we employ recurrent neural network (RNN) to propagate contextual information between utterances in a video clip, which significantly improves the classification results and outperforms the state of the art by a significant margin of 1–2% for all the modality combinations.

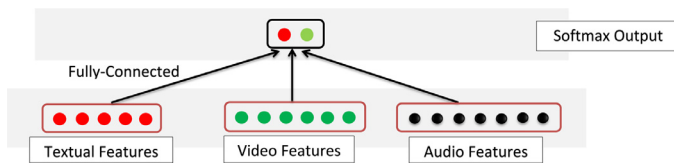
In our method, we first obtain unimodal features for each utterance for all three modalities. Then, using RNN we extract context-aware utterance features. Thus, we transform the context-aware utterance vectors to the vectors of the same dimensionality. We assume that these transformed vectors contain abstract features representing the attributes relevant to sentiment classification. Next, we compare and combine each bimodal combination of these abstract features using fully-connected layers. This yields fused bimodal feature vectors. Similarly to the unimodal case, we use RNN to generate context-aware features. Finally, we combine these bimodal vectors into a trimodal vector using, again, fully-connected layers and use a RNN to pass contextual information between them. We empirically show that the feature vectors obtained in this manner are more useful for the sentiment classification task.

The implementation of our method is publicly available in the form of open-source code.¹

* Corresponding author.

E-mail address: cambria@ntu.edu.sg (E. Cambria).

¹ <http://github.com/senticnet>.



This paper is structured as follows: [Section 2](#) briefly discusses important previous work in multimodal feature fusion; [Section 3](#) describes our method in details; [Section 4](#) reports the results of our experiments and discuss their implications; finally, [Section 5](#) concludes the paper and discusses future work.

2. Related work

In recent years, sentiment analysis has become increasingly popular for processing social media data on online communities, blogs, wikis, microblogging platforms, and other online collaborative media [2]. Sentiment analysis is a branch of affective computing research [3] that aims to classify text – but sometimes also audio and video [4] – into either positive or negative – but sometimes also neutral [5]. Most of the literature is on English language but recently an increasing number of works are tackling the multilinguality issue [6], especially in booming online languages such as Chinese [7]. Sentiment analysis techniques can be broadly categorized into symbolic and sub-symbolic approaches: the former include the use of lexicons [8], ontologies [9], and semantic networks [10] to encode the polarity associated with words and multiword expressions; the latter consist of supervised [11], semi-supervised [12] and unsupervised [13] machine learning techniques that perform sentiment classification based on word co-occurrence frequencies. Among these, the most popular recently are algorithms based on deep neural networks [14] and generative adversarial networks [15].

While most works approach it as a simple categorization problem, sentiment analysis is actually a suitcase research problem [16] that requires tackling many NLP tasks, including word polarity disambiguation [17], subjectivity detection [18], personality recognition [19], microtext normalization [20], concept extraction [21], time tagging [22], and aspect extraction [23].

Sentiment analysis has raised growing interest both within the scientific community, leading to many exciting open challenges, as well as in the business world, due to the remarkable benefits to be had from financial [24] and political [25] forecasting, e-health [26] and e-tourism [27], user profiling [28] and community detection [29], manufacturing and supply chain applications [30], human communication comprehension [31] and dialogue systems [32], etc.

In the field of emotion recognition, early works by De Silva et al. [33] and Chen et al. [34] showed that fusion of audio and visual systems, creating a bimodal signal, yielded a higher accuracy than any unimodal system. Such fusion has been analyzed at both feature level [35] and decision level [36].

Although there is much work done on audio-visual fusion for emotion recognition, exploring contribution of text along with audio and visual modalities in multimodal emotion detection has been little explored. Wollmer et al. [37] and Rozgic et al. [38] fused information from audio, visual and textual modalities to extract emotion and sentiment. Metallinou et al. [39] and Eyben et al. [40] fused audio and textual modalities for emotion recognition. Both approaches relied on a feature-level fusion. Wu and Liang [41] fused audio and textual clues at decision level. Poria et al. [42] uses convolutional neural network (CNN) to extract features from the modalities and then employs multiple-kernel learning (MKL) for sentiment analysis. The current state of the art, set forth by Poria et al. [1], extracts contextual information from the surrounding utterances using long short-term memory (LSTM). Poria et al. [3] fuses different modalities with deep learning-based

Fig. 1. Utterance-level early fusion, or simple concatenation.

tools. Zadeh et al. [43] uses tensor fusion. Poria et al. [44] further extends upon the ensemble of CNN and MKL.

Unlike existing approaches, which use simple concatenation based early fusion [42,45] and non-trainable tensors based fusion [43], this work proposes a hierarchical fusion capable of learning the bimodal and trimodal correlations for data fusion using deep neural networks. The method is end-to-end and, in order to accomplish the fusion, it can be plugged into any deep neural network based multimodal sentiment analysis framework.

3. Our method

In this section, we discuss our novel methodology behind solving the sentiment classification problem. First we discuss the overview of our method and then we discuss the whole method in details, step by step.

3.1. Overview

3.1.1. Unimodal feature extraction

We extract utterance-level features for three modalities. This step is discussed in [Section 3.2](#).

3.1.2. Multimodal fusion

Problems of early fusion. The majority of the work on multimodal data use concatenation, or early fusion ([Fig. 1](#)), as their fusion strategy. The problem with this simplistic approach is that it cannot filter out and conflicting or redundant information obtained from different modalities. To address this major issue, we devise an hierarchical approach which proceeds from unimodal to bimodal vectors and then bimodal to trimodal vectors.

Bimodal fusion. We fuse the utterance feature vectors for each bimodal combination, i.e., T + V, T + A, and A + V. This step is depicted in [Fig. 2](#) and discussed in details in [Section 3.4.1](#). We use the penultimate layer for [Fig. 2](#) as bimodal features.

Trimodal fusion. We fuse the three bimodal features to obtain trimodal feature as depicted in [Fig. 3](#).² This step is discussed in details in [Section 3.4.2](#).

Addition of context. We also improve the quality of feature vectors (both unimodal and multimodal) by incorporating information from surrounding utterances using RNN. We model the context using gated recurrent unit (GRU) as depicted in [Fig. 4](#). The details of context modeling is discussed in [Section 3.3](#) and the following subsections.

Classification. We classify the feature vectors using a softmax layer.

3.2. Unimodal feature extraction

In this section, we discuss the method of feature extraction for three different modalities: audio, video, and text.

3.2.1. Textual feature extraction

The textual data is obtained from the transcripts of the videos. We

² Figure adapted from [51] with permission.

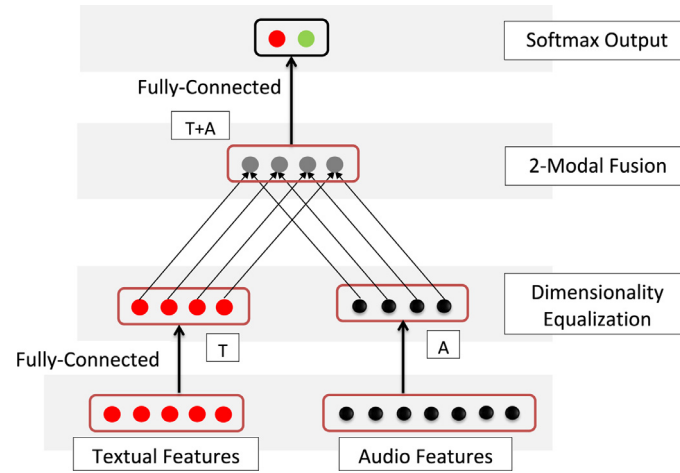


Fig. 2. Utterance-level bimodal fusion.

apply a deep Convolutional Neural Networks (CNN) [46] on each utterance to extract textual features. Each utterance in the text is represented as an array of pre-trained 300-dimensional word2vec vectors [47]. Further, the utterances are truncated or padded with null vectors to have exactly 50 words.

Next, these utterances as array of vectors are passed through two different convolutional layers; first layer having two filters of size 3 and 4 respectively with 50 feature maps each and the second layer has a filter of size 2 with 100 feature maps. Each convolutional layer is followed by a max-pooling layer with window 2×2 .

The output of the second max-pooling layer is fed to a fully-connected layer with 500 neurons with a rectified linear unit (ReLU) [48] activation, followed by softmax output. The output of the penultimate fully-connected layer is used as the textual feature. The translation of convolution filter over makes the CNN learn abstract features and with each subsequent layer the context of the features expands further.

3.2.2. Audio feature extraction

The audio feature extraction process is performed at 30 Hz frame rate with 100 ms sliding window. We use openSMILE [49], which is capable of automatic pitch and voice intensity extraction, for audio feature extraction. Prior to feature extraction audio signals are

processed with voice intensity thresholding and voice normalization. Specifically, we use Z-standardization for voice normalization. In order to filter out audio segments without voice, we threshold voice intensity. OpenSMILE is used to perform both these steps. Using openSMILE we extract several Low Level Descriptors (LLD) (e.g., pitch, voice intensity) and various statistical functionals of them (e.g., amplitude mean, arithmetic mean, root quadratic mean, standard deviation, flatness, skewness, kurtosis, quartiles, inter-quartile ranges, and linear regression slope). “IS13-ComParE” configuration file of openSMILE is used to for our purposes. Finally, we extracted total 6392 features from each input audio segment.

3.2.3. Visual feature extraction

To extract visual features, we focus not only on feature extraction from each video frame but also try to model temporal features across frames. To achieve this, we use 3D-CNN on the video. 3D-CNNs have been successful in the past, specially in the field of object classification on 3D data [50]. Its state-of-the-art performance on such tasks motivates its use in this paper.

Let the video be called $vid \in \mathbb{R}^{3 \times f \times h \times w}$, where 3 represents the three RGB channels of an image and f, h , and w denote the cardinality, height, and width of the frames, respectively. A 3D convolutional filter,

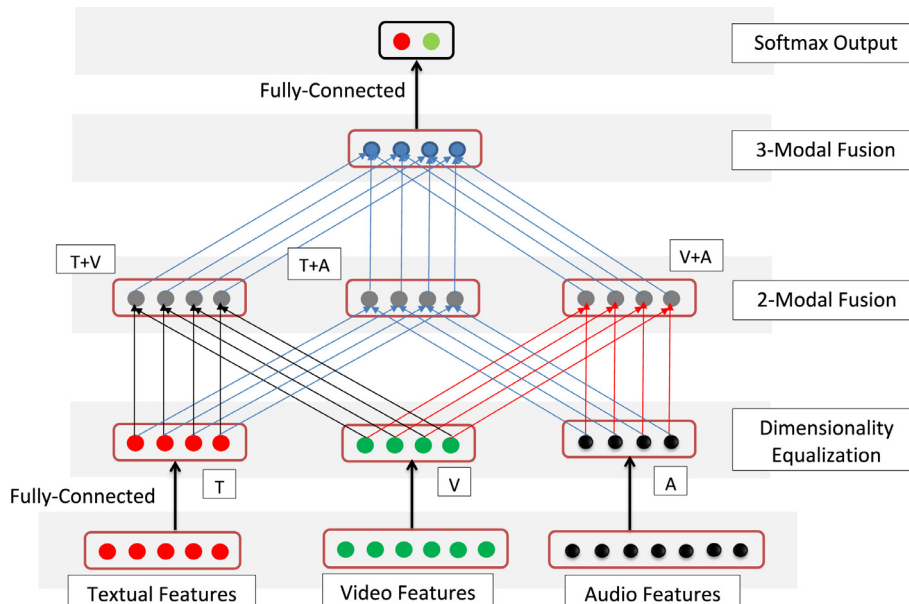


Fig. 3. Utterance-level trimodal hierarchical fusion.

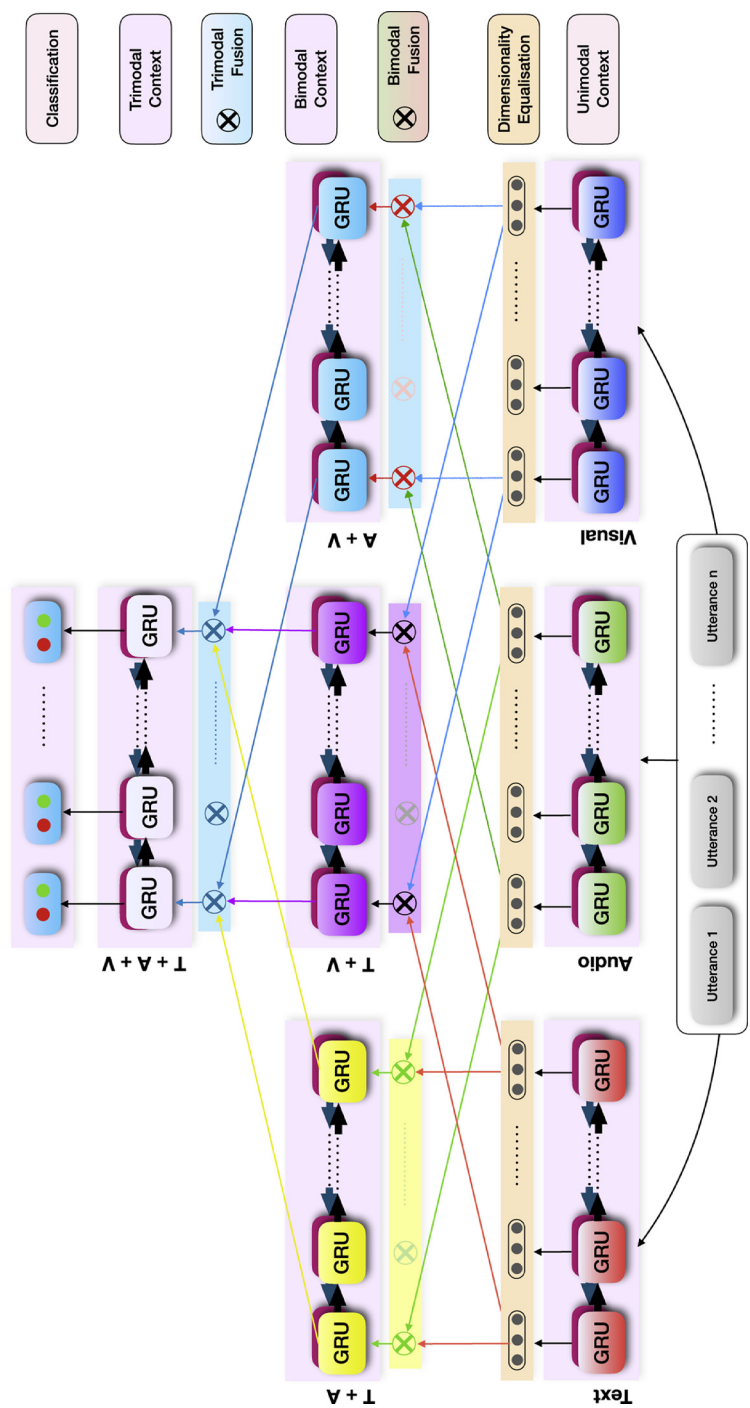


Fig. 4. Context-aware hierarchical fusion.

named $f_{lt} \in \mathbb{R}^{f_m \times 3 \times f_d \times f_h \times f_w}$, is applied to this video, where, similar to a 2D-CNN, the filter translates across the video and generates the convolution output $conv_{out} \in \mathbb{R}^{f_m \times 3 \times (f_d+1) \times (f_h+1) \times (f_w+1)}$. Here, f_m, f_d, f_h , and f_w denote number of feature maps, depth of filter, height of filter, and width of filter, respectively. Finally, we apply max-pooling operation to the $conv_{out}$, which selects the most relevant features. This operation is applied only to the last three dimensions of $conv_{out}$. This is followed by a dense layer and softmax computation. The activations of this layer is used as the overall video features for each utterance video.

In our experiments, we receive the best results with filter dimensions $f_m = 32$ and $f_d, f_h, f_w = 5$. Also, for the max-pooling, we set the window size as $3 \times 3 \times 3$ and the succeeding dense layer with 300 neurons.

3.3. Context modeling

Utterances in the videos are semantically dependent on each other. In other words, complete meaning of an utterance may be determined by taking preceding utterances into consideration. We call this the context of an utterance. Following Poria et al. [11], we use RNN, specifically GRU³ to model semantic dependency among the utterances in a video.

Let the following items represent unimodal features:

$$\begin{aligned} f_A &\in \mathbb{R}^{N \times d_A} && \text{(acoustic features),} \\ f_V &\in \mathbb{R}^{N \times d_V} && \text{(visual features),} \\ f_T &\in \mathbb{R}^{N \times d_T} && \text{(textual features),} \end{aligned}$$

where N = maximum number of utterances in a video. We pad the shorter videos with dummy utterances represented by null vectors of corresponding length. For each modality, we feed the unimodal utterance features f_m (where $m \in \{A, V, T\}$) (discussed in Section 3.2) of a video to GRU_m with output size D_m , which is defined as

$$\begin{aligned} z_m &= \sigma(f_{mt} U^{mz} + s_{m(t-1)} W^{mz}), \\ r_m &= \sigma(f_{mt} U^{mr} + s_{m(t-1)} W^{mr}), \\ h_{mt} &= \tanh(f_{mt} U^{mh} + (s_{m(t-1)} * r_m) W^{mh}), \\ F_{mt} &= \tanh(h_{mt} U^{mx} + u^{mx}), \\ s_{mt} &= (1 - z_m) * F_{mt} + z_m * s_{m(t-1)}, \end{aligned}$$

where $U^{mz} \in \mathbb{R}^{d_m \times D_m}$, $W^{mz} \in \mathbb{R}^{D_m \times D_m}$, $U^{mr} \in \mathbb{R}^{d_m \times D_m}$, $W^{mr} \in \mathbb{R}^{D_m \times D_m}$, $U^{mh} \in \mathbb{R}^{d_m \times D_m}$, $W^{mh} \in \mathbb{R}^{D_m \times D_m}$, $U^{mx} \in \mathbb{R}^{d_m \times D_m}$, $u^{mx} \in \mathbb{R}^{D_m}$, $z_m \in \mathbb{R}^{D_m}$, $r_m \in \mathbb{R}^{D_m}$, $h_{mt} \in \mathbb{R}^{D_m}$, $F_{mt} \in \mathbb{R}^{D_m}$, and $s_{mt} \in \mathbb{R}^{D_m}$. This yields hidden outputs F_{mt} as context-aware unimodal features for each modality. Hence, we define $F_m = GRU_m(f_m)$, where $F_m \in \mathbb{R}^{N \times D_m}$. Thus, the context-aware multimodal features can be defined as

$$\begin{aligned} F_A &= GRU_A(f_A), \\ F_V &= GRU_V(f_V), \\ F_T &= GRU_T(f_T). \end{aligned}$$

3.4. Multimodal fusion

In this section, we use context-aware unimodal features F_A, F_V , and F_T to a unified feature space.

The unimodal features may have different dimensions, i.e., $D_A \neq D_V \neq D_T$. Thus, we map them to the same dimension, say D (we obtained best results with $D = 400$), using fully-connected layer as follows:

$$\begin{aligned} g_A &= \tanh(F_A W_A + b_A), \\ g_V &= \tanh(F_V W_V + b_V), \\ g_T &= \tanh(F_T W_T + b_T), \end{aligned}$$

where $W_A \in \mathbb{R}^{D_A \times D}$, $b_A \in \mathbb{R}^D$, $W_V \in \mathbb{R}^{D_V \times D}$, $b_V \in \mathbb{R}^D$, $W_T \in \mathbb{R}^{D_T \times D}$, and $b_T \in \mathbb{R}^D$. We can represent the mapping for each dimension as

$$g_x = \begin{bmatrix} c_{11}^x & c_{21}^x & c_{31}^x & \cdots & c_{D1}^x \\ c_{12}^x & c_{22}^x & c_{32}^x & \cdots & c_{D2}^x \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ c_{1N}^x & c_{2N}^x & c_{3N}^x & \cdots & c_{DN}^x \end{bmatrix},$$

where $x \in \{V, A, T\}$ and c_{lt}^x are scalars for all $l = 1, 2, \dots, D$ and $t = 1, 2, \dots, N$. Also, in g_x the rows represent the utterances and the columns the feature values. We can see these values c_{lt}^x as more abstract feature values derived from fundamental feature values (which are the components of f_A, f_V , and f_T). For example, an abstract feature can be the angeriness of a speaker in a video. We can infer the degree of angeriness from visual features (f_V ; facial muscle movements), acoustic features (f_A , such as pitch and raised voice), or textual features (f_T , such as the language and choice of words). Therefore, the degree of angeriness can be represented by c_{lt}^x , where x is A, V , or T , l is some fixed integer between 1 and D , and t is some fixed integer between 1 and N .

Now, the evaluation of abstract feature values from all the modalities may not have the same merit or may even contradict each other. Hence, we need the network to make comparison among the feature values derived from different modalities to make a more refined evaluation of the degree of anger. To this end, we take each bimodal combination (which are audio–video, audio–text, and video–text) at a time and compare and combine each of their respective abstract feature values (i.e. c_{lt}^V with c_{lt}^A , c_{lt}^V with c_{lt}^T , and c_{lt}^A with c_{lt}^T) using fully-connected layers as follows:

$$i_{lt}^{VA} = \tanh(w_l^{VA} \cdot [c_{lt}^V, c_{lt}^A]^T + b_l^{VA}), \quad (1)$$

$$i_{lt}^{AT} = \tanh(w_l^{AT} \cdot [c_{lt}^A, c_{lt}^T]^T + b_l^{AT}), \quad (2)$$

$$i_{lt}^{VT} = \tanh(w_l^{VT} \cdot [c_{lt}^V, c_{lt}^T]^T + b_l^{VT}), \quad (3)$$

where $w_l^{VA} \in \mathbb{R}^2$, b_l^{VA} is scalar, $w_l^{AT} \in \mathbb{R}^2$, b_l^{AT} is scalar, $w_l^{VT} \in \mathbb{R}^2$, and b_l^{VT} is scalar, for all $l = 1, 2, \dots, D$ and $t = 1, 2, \dots, N$. We hypothesize that it will enable the network to compare the decisions from each modality against the others and help achieve a better fusion of modalities.

3.4.1. Bimodal fusion

Eqs. (1)(3) are used for bimodal fusion. The bimodal fused features for video–audio, audio–text, video–text are defined as

$$\begin{aligned} f_{VA} &= (f_{VA1}, f_{VA2}, \dots, f_{VA(N)}), \text{ where } f_{VAi} = (i_{1i}^{VA}, i_{2i}^{VA}, \dots, i_{Di}^{VA}), \\ f_{AT} &= (f_{AT1}, f_{AT2}, \dots, f_{AT(N)}), \text{ where } f_{ATi} = (i_{1i}^{AT}, i_{2i}^{AT}, \dots, i_{Di}^{AT}), \\ f_{VT} &= (f_{VT1}, f_{VT2}, \dots, f_{VT(N)}), \text{ where } f_{VTi} = (i_{1i}^{VT}, i_{2i}^{VT}, \dots, i_{Di}^{VT}). \end{aligned}$$

We further employ GRU_m (Section 3.3) ($m \in \{VA, VT, TA\}$), to incorporate contextual information among the utterances in a video with

$$\begin{aligned} F_{VA} &= (F_{VA1}, F_{VA2}, \dots, F_{VA(N)}) = GRU_{VA}(f_{VA}), \\ F_{VT} &= (F_{VT1}, F_{VT2}, \dots, F_{VT(N)}) = GRU_{VT}(f_{VT}), \\ F_{TA} &= (F_{TA1}, F_{TA2}, \dots, F_{TA(N)}) = GRU_{TA}(f_{TA}), \end{aligned}$$

where

$$\begin{aligned} F_{VAi} &= (I_{1i}^{VA}, I_{2i}^{VA}, \dots, I_{D2i}^{VA}), \\ F_{VTi} &= (I_{1i}^{AT}, I_{2i}^{AT}, \dots, I_{D2i}^{AT}), \\ F_{TAi} &= (I_{1i}^{VT}, I_{2i}^{VT}, \dots, I_{D2i}^{VT}), \end{aligned}$$

F_{VA}, F_{VT} , and F_{TA} are context-aware bimodal features represented as vectors and I_{nt}^m is scalar for $n = 1, 2, \dots, D_2$, $D_2 = 500$, $t = 1, 2, \dots, N$, and $m = VA, VT, TA$.

3.4.2. Trimodal fusion

We combine all three modalities using fully-connected layers as

³ LSTM does not perform well.

follows:

$$z_{lt} = \tanh(w_l^{AVT} \cdot [I_{lt}^{VA}, I_{lt}^{AT}, I_{lt}^{VT}]^T + b_l^{AVT}),$$

where $w_l^{AVT} \in \mathbb{R}^3$ and b_l^{AVT} is a scalar for all $l = 1, 2, \dots, D_2$ and $t = 1, 2, \dots, N$. So, we define the fused features as

$$f_{AVT} = (f_{AVT1}, f_{AVT2}, \dots, f_{AVT(N)}),$$

where $f_{AVTt} = (z_{1t}, z_{2t}, \dots, z_{D_2t})$, z_{nt} is scalar for $n = 1, 2, \dots, D_2$ and $t = 1, 2, \dots, N$.

Similarly to bimodal fusion (Section 3.4.1), after trimodal fusion we pass the fused features through GRU_{AVT} to incorporate contextual information in them, which yields

$$F_{AVT} = (F_{AVT1}, F_{AVT2}, \dots, F_{AVT(N)}) = GRU_{AVT}(f_{AVT}),$$

where $F_{AVTt} = (Z_{1t}, Z_{2t}, \dots, Z_{D_3t})$, Z_{nt} is scalar for $n = 1, 2, \dots, D_3$, $D_3 = 550$, $t = 1, 2, \dots, N$, and F_{AVT} is the context-aware trimodal feature vector.

3.5. Classification

In order to perform classification, we feed the fused features F_{mt} (where $m = AV, VT, TA$, or AVT and $t = 1, 2, \dots, N$) to a softmax layer with $C = 2$ outputs. The classifier can be described as follows:

$$\begin{aligned} \mathcal{P} &= \text{softmax}(W_{\text{softmax}} F_{mt} + b_{\text{softmax}}), \\ \hat{y} &= \underset{j}{\text{argmax}}(\mathcal{P}[j]), \end{aligned}$$

where $W_{\text{softmax}} \in \mathbb{R}^{C \times D}$, $b_{\text{softmax}} \in \mathbb{R}^C$, $\mathcal{P} \in \mathbb{R}^C$, $j = \text{class value (0 or 1)}$, and $\hat{y} = \text{estimated class value}$.

3.6. Training

We employ categorical cross-entropy as loss function (J) for training,

$$J = -\frac{1}{N} \sum_{i=1}^N \sum_{j=0}^{C-1} y_{ij} \log \mathcal{P}_i[j],$$

where $N = \text{number of samples}$, $i = \text{index of a sample}$, $j = \text{class value}$, and

$$y_{ij} = \begin{cases} 1, & \text{if expected class value of sample } i \text{ is } j \\ 0, & \text{otherwise.} \end{cases}$$

Adam [52] is used as optimizer due to its ability to adapt learning rate for each parameter individually. We train the network for 200 epochs with early stopping, where we optimize the parameter set

$$\begin{aligned} \Theta = & \bigcup_{m \in M} \left(\bigcup_{j \in \{z, r, h\}} \{U^{mj}, W^{mj}\} \cup \{U^{mx}, U^{mx}\} \right) \\ & \cup \bigcup_{m \in M_2} \bigcup_{i=1}^{D_2} \{w_i^m\} \cup \bigcup_{i=1}^{D_3} \{w_i^{AVT}\} \cup \bigcup_{m \in M_1} \{W_m, b_m\} \\ & \cup \{W_{\text{softmax}}, b_{\text{softmax}}\}, \end{aligned}$$

where $M = \{A, V, T, VA, VT, TA, AVT\}$, $M_1 = \{A, V, T\}$, and $M_2 = \{VA, VT, TA\}$. Algorithm 1 summarizes our method.⁴

4. Experiments

4.1. Dataset details

Most research works in multimodal sentiment analysis are performed on datasets where train and test splits may share certain speakers. Since, each individual has an unique way of expressing

emotions and sentiments, finding generic and person-independent features for sentiment analysis is crucial. Table 1 shows the train and test split for the datasets used.

4.1.1. CMU-MOSI

CMU-MOSI dataset [53] is rich in sentimental expressions, where 89 people review various topics in English. The videos are segmented into utterances where each utterance is annotated with scores between -3 (strongly negative) and $+3$ (strongly positive) by five annotators. We took the average of these five annotations as the sentiment polarity and considered only two classes (positive and negative). Given every individual's unique way of expressing sentiments, real world applications should be able to model generic person independent features and be robust to person variance. To this end, we perform person-independent experiments to emulate unseen conditions. Our train/test splits of the dataset are completely disjoint with respect to speakers. The train/validation set consists of the first 62 individuals in the dataset. The test set contains opinionated videos by rest of the 31 speakers. In particular, 1447 and 752 utterances are used for training and test respectively.

4.1.2. IEMOCAP

IEMOCAP [54] contains two way conversations among ten speakers, segmented into utterances. The utterances are tagged with the labels anger, happiness, sadness, neutral, excitement, frustration, fear, surprise, and other. We consider the first four ones to compare with the state of the art [1] and other works. It contains 1083 angry, 1630 happy, 1083 sad, and 1683 neutral videos. Only the videos by the first eight speakers are considered for training.

4.2. Baselines

We compare our method with the following strong baselines.

Early fusion. We extract unimodal features (Section 3.2) and simply concatenate them to produce multimodal features. Followed by support vector machine (SVM) being applied on this feature vector for the final sentiment classification.

Method from [42]. We have implemented and compared our method with the approach proposed by Poria et al. [42]. In their approach, they extracted visual features using CLM-Z, audio features using openSMILE, and textual features using CNN. MKL was then applied to the features obtained from concatenation of the unimodal features. However, they did not conduct speaker independent experiments.

In order to perform a fair comparison with [42], we employ our fusion method on the features extracted by Poria et al. [42].

Method from [1]. We have compared our method with [45], which takes advantage of contextual information obtained from the surrounding utterances. This context modeling is achieved using LSTM. We reran the experiments of Poria et al. [45] without using SVM for classification since using SVM with neural networks is usually discouraged. This provides a fair comparison with our model which does not use SVM.

Method from [43]. In [43], they proposed a trimodal fusion method based on the tensors. We have also compared our method with their. In particular, their dataset configuration was different than us so we have adapted their publicly available code⁵ and employed that on our dataset.

⁴ Implementation of this algorithm is available at <http://github.com/senticnet>.

⁵ <https://github.com/A2Zadeh/TensorFusionNetwork>.

```

1: procedure TRAINANDTESTMODEL( $U, V$ )
2:   Unimodal feature extraction:
3:   for  $i$ :  $[1, N]$  do
4:      $f_A^i \leftarrow \text{AudioFeatures}(u_i)$ 
5:      $f_V^i \leftarrow \text{VideoFeatures}(u_i)$ 
6:      $f_T^i \leftarrow \text{TextFeatures}(u_i)$ 
7:   for  $m \in \{A, V, T\}$  do
8:      $F_m = \text{GRU}_m(f_m)$ 
9:   Fusion:
10:     $g_A \leftarrow \text{MapToSpace}(F_A)$ 
11:     $g_V \leftarrow \text{MapToSpace}(F_V)$ 
12:     $g_T \leftarrow \text{MapToSpace}(F_T)$ 
13:     $f_{VA} \leftarrow \text{BimodalFusion}(g_V, g_A)$ 
14:     $f_{AT} \leftarrow \text{BimodalFusion}(g_A, g_T)$ 
15:     $f_{VT} \leftarrow \text{BimodalFusion}(g_V, g_T)$ 
16:    for  $m \in \{VA, AT, VT\}$  do
17:       $F_m = \text{GRU}_m(f_m)$ 
18:     $f_{AVT} \leftarrow \text{TrimodalFusion}(F_{VA}, F_{AT}, F_{VT})$ 
19:     $F_{AVT} = \text{GRU}_{AVT}(f_{AVT})$ 
20:    for  $i$ :  $[1, N]$  do
21:       $\hat{y}^i = \text{argmax}_j(\text{softmax}(F_{AVT}^i)[j])$ 
22:     $\text{TestModel}(V)$ 
23:  procedure MAPToSPACE( $x_z$ )
24:     $g_z \leftarrow \tanh(W_z x_z + b_z)$ 
25:    return  $g_z$ 
26:  procedure BIMODALFUSION( $g_{z_1}, g_{z_2}$ )
27:    for  $i$ :  $[1, D]$  do
28:       $f_{z_1 z_2}^i \leftarrow \tanh(w_i^{z_1 z_2} \cdot [g_{z_1}^i, g_{z_2}^i]^T + b_i^{z_1 z_2})$ 
29:       $f_{z_1 z_2} \leftarrow (f_{z_1 z_2}^1, f_{z_1 z_2}^2, \dots, f_{z_1 z_2}^D)$ 
30:    return  $f_{z_1 z_2}$ 
31:  procedure TRIMODALFUSION( $f_{z_1}, f_{z_2}, f_{z_3}$ )
32:    for  $i$ :  $[1, D]$  do
33:       $f_{z_1 z_2 z_3}^i \leftarrow \tanh(w_i^{z_1 z_2 z_3} \cdot [f_{z_1}^i, f_{z_2}^i, f_{z_3}^i]^T + b_i)$ 
34:       $f_{z_1 z_2 z_3} \leftarrow (f_{z_1 z_2 z_3}^1, f_{z_1 z_2 z_3}^2, \dots, f_{z_1 z_2 z_3}^D)$ 
35:    return  $f_{z_1 z_2 z_3}$ 
36:  procedure TESTMODEL( $V$ )
37:    Similarly to training phase,  $V$  is passed through the learnt models to get the features and classification outputs. Section 3.6 mentions the trainable parameters ( $\theta$ ).

```

▶ U = train set, V = test set
 ▶ extract baseline features
 ▶ dimensionality equalization
 ▶ bimodal fusion
 ▶ trimodal fusion
 ▶ softmax classification
 ▶ for modality z
 ▶ for modality z_1 and z_2 , where $z_1 \neq z_2$
 ▶ for modality combination z_1, z_2 , and z_3 , where $z_1 \neq z_2 \neq z_3$

Algorithm 1. Context-aware hierarchical fusion algorithm.

Table 1
Class distribution of datasets in both train and test splits.

Dataset	Train						Test					
	pos.	neg.	happy	anger	sad	neu.	pos.	neg.	happy	anger	sad	neu.
MOSI	709	738	–	–	–	–	467	285	–	–	–	–
IEMOCAP	–	–	1194	933	839	1324	–	–	433	157	238	380

pos. = positive, neg. = negative, neu. = neutral

4.3. Experimental setting

We considered two variants of experimental setup while evaluating our model.

HFusion. In this setup, we evaluated hierarchical fusion without context-aware features with CMU-MOSI dataset. We removed all the GRUs from the model described in Sections 3.4 and 3.3 forwarded utterance specific features directly to the next layer. This setup is depicted in Fig. 3.

CHFusion. This setup is exactly as the model described in Section 3.

4.4. Results and discussion

We discuss the results for the different experimental settings discussed in Section 4.3.

4.4.1. Hierarchical fusion (HFusion)

The results of our experiments are presented in Table 2. We evaluated this setup with CMU-MOSI dataset (Section 4.1.1) and two feature sets: the feature set used in [42] and the set of unimodal features discussed in Section 3.2.

Our model outperformed [42], which employed MKL, for all bimodal and trimodal scenarios by a margin of 1–1.8%. This leads us to present two observations. Firstly, the features used in [42] are inferior to the features extracted in our approach. Second, our hierarchical fusion method is better than their fusion method.

It is already established in the literature [42,55] that multimodal analysis outperforms unimodal analysis. We also observe the same trend in our experiments where trimodal and bimodal classifiers outperform unimodal classifiers. The textual modality performed best among others with a higher unimodal classification accuracy of 75%. Although other modalities contribute to improve the performance of multimodal classifiers, that contribution is little in compare to the textual modality.

On the other hand, we compared our model with early fusion

Table 2

Comparison in terms of accuracy of Hierarchical Fusion (HFusion) with other fusion methods for CMU-MOSI dataset; bold font signifies best accuracy for the corresponding feature set and modality or modalities, where T stands for text, V for video, and A for audio. $SOTA^1$ = Poria et al. [42], $SOTA^2$ = Zadeh et al. [43].

Modality combination	[42] feature set			Our feature set		
	$SOTA^1$	$SOTA^2$	HFusion	Early fusion	$SOTA^2$	HFusion
T	N/A			75.0%		
V	N/A			55.3%		
A	N/A			56.9%		
T + V	73.2%	73.8%	74.4%	77.1%	77.4%	77.8%
T + A	73.2%	73.5%	74.2%	77.1%	76.3%	77.3%
A + V	55.7%	56.2%	57.5%	56.5%	56.1%	56.8%
A + V + T	73.5%	71.2%	74.6%	77.0%	77.3%	77.9%

(Section 4.2) for aforementioned feature sets (Section 3.2). Our fusion mechanism consistently outperforms early fusion for all combination of modalities. This supports our hypothesis that our hierarchical fusion method captures the inter-relation among the modalities and produce better performance vector than early fusion. Text is the strongest individual modality, and we observe that the text modality paired with remaining two modalities results in consistent performance improvement.

Overall, the results give a strong indication that the comparison among the abstract feature values dampens the effect of less important modalities, which was our hypothesis. For example, we can notice that for early fusion T + V and T + A both yield the same performance. However, with our method text with video performs better than text with audio, which is more aligned with our expectations, since facial muscle movements usually carry more emotional nuances than voice.

In particular, we observe that our model outperformed all the strong baselines mentioned above. The method by Poria et al. [42] is only able to fuse using concatenation. Our proposed method outperformed their approach by a significant margin; thanks to the power of hierarchical fusion which proves the capability of our method in modeling bimodal and trimodal correlations. However on the other hand, the method by Zadeh et al. [43] is capable of fusing the modalities using a tensor. Interestingly our method also outperformed them and we think the reason is the capability of bimodal fusion and use that for trimodal fusion. Tensor fusion network is incapable to learn the weights of the bimodal and trimodal correlations in the fusion. Tensor Fusion is mathematically formed by an outer product, it has no learn-able parameters. Wherein our method learns the weights automatically using a neural network (Eq. (1)–(3)).

4.4.2. Context-aware hierarchical fusion (CHFusion)

The results of this experiment are shown in Table 3. This setting fully utilizes the model described in Section 3. We applied this experimental setting for two datasets, namely CMU-MOSI (Section 4.1.1) and IEMOCAP (Section 4.1.2). We used the feature set discussed in Section 3.2, which was also used by Poria et al. [1]. As expected our method outperformed the simple early fusion based fusion by Poria et al. [42], tensor fusion by Zadeh et al. [43]. The method by Poria et al. [1] used a scheme to learn contextual features from the surrounding features. However, as a method of fusion they adapted simple concatenation based fusion method by Poria et al. [42]. As discussed in Section 3.3, we employed their contextual feature extraction framework and integrated our proposed fusion method to that. This has helped us to outperform Poria et al. [1] by significant margin thanks to the hierarchical fusion (HFusion).

CMU-MOSI. We achieve 1–2% performance improvement over the state of the art [1] for all the modality combinations having textual component. For A + V modality combination we achieve better but similar performance to the state of the art. We suspect that it is due to both audio and video modality being significantly less informative than textual modality. It is evident from the unimodal performance where we observe that textual modality on its own performs around 21% better than both audio and video modality. Also, audio and video modality performs close to majority baseline. On the other hand, it is

Table 3

Comparison of Context-Aware Hierarchical Fusion (CHFusion) in terms of accuracy (CHFusion_{acc}) and f-score (for IEMOCAP: CHFusion_{fsc}) with the state of the art for CMU-MOSI and IEMOCAP dataset; bold font signifies best accuracy for the corresponding dataset and modality or modalities, where T stands text, V for video, A for audio. $SOTA^1$ = Poria et al. [42], $SOTA^2$ = Zadeh et al. [43]. CHFusion_{acc} and CHFusion_{fsc} are the accuracy and f-score of CHFusion respectively.

Modality	CMU-MOSI			IEMOCAP			
	$SOTA^1$	$SOTA^2$	CHFusion _{acc}	$SOTA^1$	$SOTA^2$	CHFusion _{acc}	CHFusion _{fsc}
T	76.5%			73.6%			–
V	54.9%			53.3%			–
A	55.3%			57.1%			–
T + V	77.8%	77.1%	79.3%	74.1%	73.7%	75.9%	75.6%
T + A	77.3%	77.0%	79.1%	73.7%	71.1%	76.1%	76.0%
A + V	57.9%	56.5%	58.8%	68.4%	67.4%	69.5%	69.6%
A + V + T	78.7%	77.2%	80.0%	74.1%	73.6%	76.5%	76.8%

Table 4

Class-wise accuracy and f-score for IEMOCAP dataset for trimodal scenario.

Metrics	Classes			
	Happy	Sad	Neutral	Anger
Accuracy	74.3	75.6	78.4	79.6
F-Score	81.4	77.0	71.2	77.6

important to notice that with all modalities combined we achieve about 3.5% higher accuracy than text alone.

For example, consider the following utterance: *so overall new moon even with the bigger better budgets huh it was still too long*. The speaker discusses her opinion on the movie Twilight New Moon. Textually the utterance is abundant with positive words however audio and video comprises of a frown which is observed by the hierarchical fusion based model.

IEMOCAP. As the IEMOCAP dataset contains four distinct emotion categories, in the last layer of the network we used a softmax classifier whose output dimension is set to 4. In order to perform classification on IEMOCAP dataset we feed the fused features F_{mt} (where $m = AV, VT, TA$, or AVT and $t = 1, 2, \dots, N$) to a softmax layer with $C = 4$ outputs. The classifier can be described as follows:

$$\mathcal{P} = \text{softmax}(W_{\text{softmax}} F_{mt} + b_{\text{softmax}}),$$

$$\hat{y} = \underset{j}{\text{argmax}}(\mathcal{P}[j]),$$

where $W_{\text{softmax}} \in \mathbb{R}^{4 \times D}$, $b_{\text{softmax}} \in \mathbb{R}^4$, $\mathcal{P} \in \mathbb{R}^4$, j = class value (0 or 1 or 2 or 3), and \hat{y} = estimated class value.

Here as well, we achieve performance improvement consistent with CMU-MOSI. This method performs 1–2.4% better than the state of the art for all the modality combinations. Also, trimodal accuracy is 3% higher than the same for textual modality. Since, IEMOCAP dataset imbalanced, we also present the f-score for each modality combination for a better evaluation. One key observation for IEMOCAP dataset is that its A + V modality combination performs significantly better than the same of CMU-MOSI dataset. We think that this is due to the audio and video modality of IEMOCAP being richer than the same of CMU-MOSI. The performance difference with another strong baseline [43] is even more ranging from 2.1% to 3% on CMU-MOSI dataset and 2.2–5% on IEMOCAP dataset. This again confirms the superiority of the hierarchical fusion in compare to Zadeh et al. [43]. We think this is mainly because of learning the weights of bimodal and trimodal correlation (representing the degree of correlations) calculations at the time of fusion while Tensor Fusion Network (TFN) just relies on the non-trainable outer product of tensors to model such correlations for fusion. Additionally, we present class-wise accuracy and f-score for IEMOCAP for trimodal (A + V + T) scenario in Table 4.

4.4.3. HFusion vs. CHFusion

We compare HFusion and CHFusion models over CMU-MOSI dataset. We observe that CHFusion performs 1–2% better than HFusion model for all the modality combinations. This performance boost is achieved by the inclusion of utterance-level contextual information in HFusion model by adding GRUs in different levels of fusion hierarchy.

5. Conclusion

Multimodal fusion strategy is an important issue in multimodal sentiment analysis. However, little work has been done so far in this direction. In this paper, we have presented a novel and comprehensive fusion strategy. Our method outperforms the widely used early fusion on both datasets typically used to test multimodal sentiment analysis methods. Moreover, with the addition of context modeling with GRU, our method outperforms the state of the art in multimodal sentiment analysis and emotion detection by significant margin.

In our future work, we plan to improve the quality of unimodal features, especially textual features, which will further improve the accuracy of classification. We will also experiment with more sophisticated network architectures.

Acknowledgment

The work was partially supported by the Instituto Politécnico Nacional via grant SIP 20172008 to A. Gelbukh.

References

- [1] S. Poria, E. Cambria, D. Hazarika, N. Mazumder, A. Zadeh, L.-P. Morency, Context-dependent sentiment analysis in user-generated videos, *ACL*, (2017), pp. 873–883.
- [2] E. Cambria, Affective computing and sentiment analysis, *IEEE Intell. Syst.* 31 (2) (2016) 102–107.
- [3] S. Poria, E. Cambria, R. Bajpai, A. Hussain, A review of affective computing: From unimodal analysis to multimodal fusion, *Inf. Fusion* 37 (2017) 98–125.
- [4] D. Hazarika, S. Poria, A. Zadeh, E. Cambria, L.-P. Morency, R. Zimmermann, Conversational memory network for emotion recognition in dyadic dialogue videos, *NAACL*, (2018), pp. 2122–2132.
- [5] I. Chaturvedi, E. Cambria, R. Welsch, F. Herrera, Distinguishing between facts and opinions for sentiment analysis: Survey and challenges, *Inf. Fusion* 44 (2018) 65–77.
- [6] S.L. Lo, E. Cambria, R. Chiong, D. Cornforth, Multilingual sentiment analysis: From formal to informal and scarce resource languages, *Artif. Intell. Rev.* 48 (4) (2017) 499–527.
- [7] H. Peng, Y. Ma, Y. Li, E. Cambria, Learning multi-grained aspect target sequence for chinese sentiment analysis, *Knowl. Based Syst.* 148 (2018) 167–176.
- [8] A. Bandhakavi, N. Wiratunga, S. Massie, P. Deepak, Lexicon generation for emotion analysis of text, *IEEE Intell. Syst.* 32 (1) (2017) 102–108.
- [9] M. Dragoni, S. Poria, E. Cambria, OntoSentNet: A commonsense ontology for sentiment analysis, *IEEE Intell. Syst.* 33 (3) (2018) 77–85.
- [10] E. Cambria, S. Poria, D. Hazarika, K. Kwok, SenticNet 5: Discovering conceptual primitives for sentiment analysis by means of context embeddings, *AAAI*, (2018), pp. 1795–1802.
- [11] L. Oneto, F. Bisio, E. Cambria, D. Anguita, Statistical learning theory and ELM for big social data analysis, *IEEE Comput. Intell. Mag.* 11 (3) (2016) 45–55.
- [12] A. Hussain, E. Cambria, Semi-supervised learning for big social data analysis,

- Neurocomputing 275 (2018) 1662–1673.
- [13] Y. Li, Q. Pan, T. Yang, S.H. Wang, J.L. Tang, E. Cambria, Learning word representations for sentiment analysis, *Cognit. Comput.* 9 (6) (2017) 843–851.
 - [14] T. Young, D. Hazarika, S. Poria, E. Cambria, Recent trends in deep learning based natural language processing, *IEEE Comput. Intell. Mag.* 13 (3) (2018) 55–75.
 - [15] Y. Li, Q. Pan, S. Wang, T. Yang, E. Cambria, A generative model for category text generation, *Inf. Sci.* 450 (2018) 301–315.
 - [16] E. Cambria, S. Poria, A. Gelbukh, M. Thelwall, Sentiment analysis is a big suitcase, *IEEE Intell. Syst.* 32 (6) (2017) 74–80.
 - [17] Y. Xia, E. Cambria, A. Hussain, H. Zhao, Word polarity disambiguation using bayesian model and opinion-level features, *Cognit. Comput.* 7 (3) (2015) 369–380.
 - [18] I. Chaturvedi, E. Ragusa, P. Gastaldo, R. Zunino, E. Cambria, Bayesian network based extreme learning machine for subjectivity detection, *J. Franklin Inst.* 355 (4) (2018) 1780–1797.
 - [19] N. Majumder, S. Poria, A. Gelbukh, E. Cambria, Deep learning-based document modeling for personality detection from text, *IEEE Intell. Syst.* 32 (2) (2017) 74–79.
 - [20] R. Satapathy, C. Guerreiro, I. Chaturvedi, E. Cambria, Phonetic-based microtext normalization for twitter sentiment analysis, *ICDM*, (2017), pp. 407–413.
 - [21] D. Rajagopal, E. Cambria, D. Olshe, K. Kwok, A graph-based approach to commonsense concept extraction and semantic similarity detection, *WWW*, (2013), pp. 565–570.
 - [22] X. Zhong, A. Sun, E. Cambria, Time expression analysis and recognition using syntactic token types and general heuristic rules, *ACL*, (2017), pp. 420–429.
 - [23] Y. Ma, H. Peng, E. Cambria, Targeted aspect-based sentiment analysis via embedding commonsense knowledge into an attentive LSTM, *AAAI*, (2018), pp. 5876–5883.
 - [24] F. Xing, E. Cambria, R. Welsch, Natural language based financial forecasting: A survey, *Artif. Intell. Rev.* 50 (1) (2018) 49–73.
 - [25] M. Ebrahimi, A. Hossein, A. Sheth, Challenges of sentiment analysis for dynamic events, *IEEE Intell. Syst.* 32 (5) (2017) 70–75.
 - [26] E. Cambria, A. Hussain, T. Durrani, C. Havasi, C. Eckl, J. Munro, Sentic computing for patient centered application, *IEEE ICSP*, (2010), pp. 1279–1282.
 - [27] A. Valdivia, V. Luzon, F. Herrera, Sentiment analysis in tripadvisor, *IEEE Intell. Syst.* 32 (4) (2017) 72–77.
 - [28] R. Mihalcea, A. Garimella, What men say, what women hear: Finding gender-specific meaning shades, *IEEE Intell. Syst.* 31 (4) (2016) 62–67.
 - [29] S. Cavallari, V. Zheng, H. Cai, K. Chang, E. Cambria, Learning community embedding with community detection and node embedding on graphs, *CIKM*, (2017), pp. 377–386.
 - [30] C. Xu, E. Cambria, P.S. Tan, Adaptive two-stage feature selection for sentiment classification, *IEEE SMC*, (2017), pp. 1238–1243.
 - [31] A. Zadeh, P.P. Liang, S. Poria, P. Vij, E. Cambria, L.-P. Morency, Multi-attention recurrent network for human communication comprehension, *AAAI*, (2018), pp. 5642–5649.
 - [32] T. Young, E. Cambria, I. Chaturvedi, H. Zhou, S. Biswas, M. Huang, Augmenting end-to-end dialogue systems with commonsense knowledge, *AAAI*, (2018), pp. 4970–4977.
 - [33] L.C. De Silva, T. Miyasato, R. Nakatsu, Facial emotion recognition using multimodal information, *Proceedings of ICICS*, vol. 1, IEEE, 1997, pp. 397–401.
 - [34] L.S. Chen, T.S. Huang, T. Miyasato, R. Nakatsu, Multimodal human emotion/expressions recognition, *Proceedings of the Third IEEE International Conference on Automatic Face and Gesture Recognition*, IEEE, 1998, pp. 366–371.
 - [35] L. Kessous, G. Castellano, G. Caridakis, Multimodal emotion recognition in speech-based interaction using facial expression, body gesture and acoustic analysis, *J. Multimodal User Interfaces* 3 (1–2) (2010) 33–48.
 - [36] B. Schuller, Recognizing affect from linguistic information in 3d continuous space, *IEEE Trans. Affect. Comput.* 2 (4) (2011) 192–205.
 - [37] M. Wollmer, F. Weninger, T. Knaup, B. Schuller, C. Sun, K. Sagae, L.-P. Morency, Youtube movie reviews: Sentiment analysis in an audio-visual context, *IEEE Intell. Syst.* 28 (3) (2013) 46–53.
 - [38] V. Rozic, S. Ananthakrishnan, S. Saleem, R. Kumar, R. Prasad, Ensemble of SVM trees for multimodal emotion recognition, *Signal & Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2012 Asia-Pacific, IEEE, 2012, pp. 1–4.
 - [39] A. Metallinou, S. Lee, S. Narayanan, Audio-visual emotion recognition using gaussian mixture models for face and voice, *Tenth IEEE International Symposium on ISM 2008*, IEEE, 2008, pp. 250–257.
 - [40] F. Eyben, M. Wöllmer, A. Graves, B. Schuller, E. Douglas-Cowie, R. Cowie, On-line emotion recognition in a 3-d activation-valence-time continuum using acoustic and linguistic cues, *J. Multimodal User Interfaces* 3 (1–2) (2010) 7–19.
 - [41] C.-H. Wu, W.-B. Liang, Emotion recognition of affective speech based on multiple classifiers using acoustic-prosodic information and semantic labels, *IEEE Trans. Affect. Comput.* 2 (1) (2011) 10–21.
 - [42] S. Poria, I. Chaturvedi, E. Cambria, A. Hussain, Convolutional MKL based multimodal emotion recognition and sentiment analysis, *ICDM*, (2016), pp. 439–448.
 - [43] A. Zadeh, M. Chen, S. Poria, E. Cambria, L.-P. Morency, Tensor fusion network for multimodal sentiment analysis, *EMNLP*, (2017), pp. 1114–1125.
 - [44] S. Poria, H. Peng, A. Hussain, N. Howard, E. Cambria, Ensemble application of convolutional neural networks and multiple kernel learning for multimodal sentiment analysis, *Neurocomputing* 261 (2017) 217–230.
 - [45] S. Poria, E. Cambria, A. Gelbukh, Deep convolutional neural network textual features and multiple kernel learning for utterance-level multimodal sentiment analysis, *EMNLP*, (2015), pp. 2539–2544.
 - [46] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, L. Fei-Fei, Large-scale video classification with convolutional neural networks, *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, (2014), pp. 1725–1732.
 - [47] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, *arXiv:1301.3781* (2013).
 - [48] V. Teh, G.E. Hinton, Rate-coded restricted boltzmann machines for face recognition, in: T. Leen, T. Dietterich, V. Tresp (Eds.), *Advances in Neural Information Processing System*, vol. 13, 2001, pp. 908–914.
 - [49] F. Eyben, M. Wöllmer, B. Schuller, Opensmile: the munich versatile and fast open-source audio feature extractor, *Proceedings of the 18th ACM International Conference on Multimedia*, ACM, 2010, pp. 1459–1462.
 - [50] S. Ji, W. Xu, M. Yang, K. Yu, 3D convolutional neural networks for human action recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (1) (2013) 221–231.
 - [51] N. Majumder, Multimodal Sentiment Analysis in Social Media using Deep Learning with Convolutional Neural Networks, *CIC*, Instituto Politécnico Nacional, 2017.
 - [52] D.P. Kingma, J. Ba, Adam: A method for stochastic optimization, *CoRR abs/1412.6980* (2014).
 - [53] A. Zadeh, R. Zellers, E. Pincus, L.-P. Morency, Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages, *IEEE Intell. Syst.* 31 (6) (2016) 82–88.
 - [54] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J.N. Chang, S. Lee, S.S. Narayanan, Iemocap: Interactive emotional dyadic motion capture database, *Lang. Resour. Eval.* 42 (4) (2008) 335–359.
 - [55] V. Pérez-Rosas, R. Mihalcea, L.-P. Morency, Utterance-level multimodal sentiment analysis, *ACL* (1), (2013), pp. 973–982.