

A Framework On Multi-Modal Sentiment Analysis Through Image Content

Ms. Pooja Morey¹, Prof. Y. B. Jadhao²

¹M.E Computer Engineering , Student, Padm. Dr. V.B.Kolte COE, Malkapur, SGBAU, INDIA

²M.E Computer Engineering , Faculty, Padm. Dr. V.B.Kolte COE, Malkapur, SGBAU, INDIA

¹poojamorey@gmail.com

²ybjadhao@gmail.com

Abstract

In latest years, with the recognition of social media, customers are an increasing number of eager to categorical their emotions and opinions in the structure of pics and text, which makes multimodal records with textual content and photos the content kind with the most growth. Most of the statistics posted via customers on social media has apparent sentimental aspects, and multimodal sentiment evaluation has end up an necessary lookup field. Previous research on multimodal sentiment evaluation have principally centered on extracting textual content and photograph aspects one by one and then combining them for sentiment classification. These research frequently bypass the interplay between textual content and images. Therefore, this challenge proposes a new multimodal sentiment evaluation model. The mannequin first eliminates noise interference in textual statistics and extracts greater necessary photo features. Then, in the feature-fusion section based totally on the interest mechanism, the textual content and pics study the interior points from every different via symmetry. Then the fusion elements are utilized to sentiment classification tasks. The experimental effects on two frequent multimodal sentiment datasets reveal the effectiveness of the proposed model.

Keywords: Data, Sentiment, Analysis, Classification, Multimodal

1. Introduction

In recent years, with the reputation of social media, customers are an increasing number of eager to specific their emotions and opinions in the structure of photographs and text, which makes multimodal records with textual content and photos the content kind with the most growth. Most of the data posted through customers on social media has apparent sentimental aspects, and multimodal sentiment evaluation has come to be an vital lookup field. Previous research on multimodal sentiment evaluation have exceptionally targeted on extracting textual content and picture elements one by one and then combining them for sentiment classification. These researches frequently omit the interplay between textual content and images. Therefore, this venture proposes a new multimodal sentiment evaluation model. The mannequin first eliminates noise interference in textual records and extracts extra vital photo features. Then, in the feature-fusion section based totally on the interest mechanism, the textual content and photographs research the inner facets from every different thru symmetry. Then the fusion aspects are utilized to sentiment classification tasks. The experimental outcomes on two frequent multimodal sentiment datasets exhibit the effectiveness of the proposed model.

The aim of photo classification is to figure out whether or not an photograph belongs to a positive class or not. Different sorts of classes have been regarded in the literature, e.g. described with the aid of presence of sure objects, such as vehicles or bicycles, or described in phrases of scene types, such as city, coast, mountain, etc. To remedy this problem, a binary classifier can be realized from a series of pix manually labeled to belong to the class or not. Increasing the extent and range of hand-labeled photos improves Tags: desert, nature, landscape, sky Tags: rose, red Labels: clouds, plant life, sky, tree Labels: flower, plant existence Tags: India Tags: aviation, airplane, airport Labels: cow Labels: aero plane. This motivates our pastime in the use of different sources of data that can resource the getting to know method the use of a restrained quantity of labeled images. With the growing reputation of social media, humans are increasingly more eager

to specific their views or opinions on social media platforms. In social media, lots of thousands and thousands of records archives are generated each day. A giant extent of records is in the shape of textual content and photograph combinations, which represent a large quantity of multimodal data. Rich sentimental facts exist in the multimodal data.

2. Literature Survey

Recognition methods based on deep neural networks have shown many advantages in terms of learning ability, high variability, and generalization. However, efficient algorithms still present several limitations when real-time operation is required, as well as in an unconstrained environments because it requires achieving high accuracy and computational efficiency. Therefore, face recognition still represents an important challenge in real-time applications, and it is an active research field in the context of computer vision, deep learning, real-time systems, etc. The complexity of face recognition systems depends on the interaction of several less complex sub-systems jointly operating to solve more complex tasks; in particular, we can generalize two fundamental operations involved in the facial recognition tasks: face detection and face recognition. A face recognition system is limited in minor conditions and required to detect faces in images (or videos) regardless of the facial object appearance. Secondly, face images are then processed; subsequently, face features are extracted with a feature extractor. Finally, the system compares the extracted features with the enrolled faces to make face matching. This is why faster algorithms can greatly benefit the system performance and achieve high recognition rates. [1]

Over the past few decades, interest in theories and algorithms for face recognition has been growing rapidly. Video surveillance, criminal identification, building access control, and unmanned and autonomous vehicles are just a few examples of concrete applications that are gaining attraction among industries. Various techniques are being developed including local, holistic, and hybrid approaches, which provide a face image description using only a few face image features or the whole facial features. The main contribution of this survey is to review some well-known techniques for each approach and to give the taxonomy of their categories. In the paper, a detailed comparison between these techniques is exposed by listing the advantages and the disadvantages of their schemes in terms of robustness, accuracy, complexity, and discrimination. One interesting feature mentioned in the paper is about the database used for face recognition. An overview of the most commonly used databases, including those of supervised and unsupervised learning, is given. Numerical results of the most interesting techniques are given along with the context of experiments and challenges handled by these techniques. Finally, a solid discussion is given in the paper [2]

Face recognition (FR) is an extensively studied topic in computer vision. Among the existing technologies of human biometrics, face recognition is the most widely used one in real-world applications. With the great advance of deep convolutional neural networks (DCNNs), the deep learning based methods have achieved significant improvements on various computer vision tasks, including face recognition. In this survey, we focus on 2D image based end-to-end deep face recognition which takes the general images or video frames as input, and extracts the deep feature of each face as output. We provide a comprehensive review of the recent advances of the elements of end-to-end deep face recognition. Specifically, an end-to-end deep face recognition system is composed of three key elements: face detection, face alignment, and face representation. In the following, we give a brief introduction of each element. In the face representation stage, the discriminative features are extracted from the aligned face images for recognition. This is the final and core step of face recognition. In early studies, many approaches calculate the face representation by projecting face images into low-dimensional subspace, such as Eigenfaces and Fisherfaces. Later on, handcrafted local descriptors based methods prevail in this area. For a detailed review of these traditional methods. [3]

With the rapid growth in multimedia contents, among such content face recognition has got much attention especially in past few years. Face as an object consists of distinct features for detection; therefore, it remains most challenging research area for scholars in the field of computer vision and image processing. In this survey paper, we have tried to address most endeavoring face features such as pose invariance, aging, illuminations and partial occlusion. They are considered to be indispensable factors in face recognition system when realized over facial images. This paper also studies state of the art face detection techniques, approaches, viz. Eigen face, Artificial Neural Networks (ANN), Support Vector Machines (SVM), Principal Component Analysis (PCA), Independent Component Analysis (ICA), Gabor Wavelets, Elastic Bunch Graph Matching, 3D morphable Model and Hidden Markov Models. In addition to the aforementioned works, we have mentioned different testing face databases which include AT & T (ORL), AR, FERET, LFW, YTF, and Yale, respectively for results analysis. [4]

Face recognition is a relatively mature technology, which has some applications in many aspects, and now there are many networks studying it, which has indeed brought a lot of convenience to mankind in all aspects. This paper proposes a new face recognition technology. First, a new GoogLeNet-M network is proposed, which improves network performance on the basis of streamlining the network. Secondly, regularization and migration learning methods are added to improve accuracy. The experimental results show that the GoogLeNet-M network with regularization using migration learning technology has the best

performance, with a recall rate of 0.97 and an accuracy of 0.98. Finally, it is concluded that the performance of the GoogLeNet-M network is better than other networks on the dataset, and the migration learning method and regularization help to improve the network performance. This has been in the era of big data, which has brought about an explosive increase in the amount of information, and in some access control and other aspects, people often use biometrics for identity authentication for a reason because people's faces or fingerprints are unique. In this regard, face recognition is the main recognition method, which brings great convenience to people's life. It mainly uses optical imaging of human faces to perceive and recognize people. [5]

Sentiment analysis uses information retrieval and computational linguistics. Sentiment analysis has advantages in various forms such as in marketing or for business purposes. In marketing, it is used to notice about the favorable or negative points about their new product which helps to determine how successful the new product is. A specific view or notion can be depicted as ideas prompted, opinions, judgements or coloured by emotions or emotions. In Computational Linguistics, the core is on feelings instead of sentiments, opinions or perceptions. The terms „opinion“s and „sentiment“s are frequently availed substitutable. In general, the information of a text is divided into two categories. 1. Based on text 2. Based on persuasion. Whereas actualities or facts are observational utterances about events, entities and their opinions, characteristics are particular utterances that depict opinions of people, events and their properties, feelings towards entities or appraisals. A persuasion can be depicted by the following four terms: Sentiment, Claim, Holder and Topic . The Holder affirms a fact about a Topic, and frequently relates a persuasion, such as 'bad' or 'good', with the affirm. It depicts a persuasion as an implicit or explicit aspect in text of the holder's negative, positive or neutral notice into the requirement about the topic. Sentiment analysis suits with computational operation of persuasion, sentiment, opinion and individuality in text. The document inception is likely in the pattern of unstructured data. [6]

3. Existing System and technologies

Traditionally the document classification was performed on the topic basis but later research started working on opinion basis. Following machine learning methods Naive Bayes, Maximum Entropy Classification (MEC), and Support

Vector Machine (SVM) are used for sentiment analysis. The conventional method of document classification based on topic is tried out for sentiment analysis. The major two classes are considered i.e. positive and negative and classify the reviews according to that. In [5], Naïve Bayes is best suitable for textual classification, clustering for consumer services and Support Vector Machine for biological reading and interpretation. The four methods discussed in the paper are actually applicable in different areas like clustering is applied in reviews and Support Vector Machine (SVM) techniques is applied in biological reviews & analysis. Though field of opinion mining is latest technology, but still it provides diverse methods available to provide a way to implement these methods.

There is a large literature on semi-supervised learning techniques. For sake of brevity, we discuss only two important paradigms, and we refer to [5] for a recent book on the subject. When using generative models for semi-supervised learning a straightforward approach is to treat the class label of unlabeled data as a missing variable, see e.g. [1, 15]. The class conditional models over the features can then be iteratively estimated using the EM algorithm. In each iteration the current model is used to estimate the class label of unlabeled data, and then the class conditional models are updated given the current label estimates. This idea can be extended to our setting where we have variables that are only observed for the training data. The idea is to jointly predict the class label and the missing text features for the test-data, and then marginalize over the unobserved text features. These methods are known to work well in cases where the model fits the data distribution, but can be detrimental in cases where the model has a poor fit. Current state-of-the-art image classification methods are discriminative ones that do not estimate the class conditional density models, but directly estimate a decision function to separate the classes. However, using discriminative classifiers, the EM method of estimating the missing class labels used for generative models does not apply: the EM iterations immediately terminate at the initial classifier. Co-training [4] is a semi-supervised learning technique that does apply to discriminative classifiers, and is designed for settings like ours where the data is described using several different feature sets. The idea is to learn a separate classifier using each feature set, and to iteratively add training examples for each classifier based on the output of the other classifier. In particular, in each iteration the examples that are most confidently classified with the first classifier are added as labeled examples to the training set of the second classifier, and vice-versa.

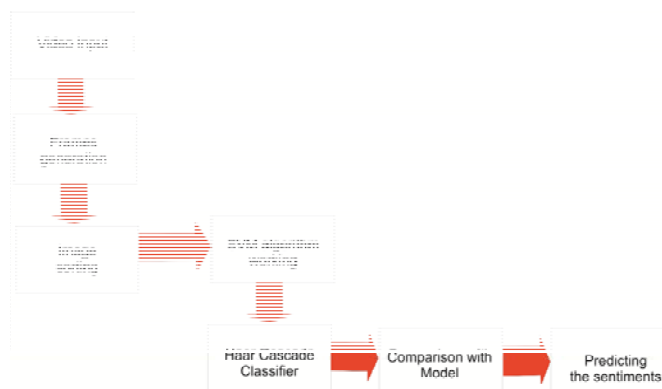


Fig. 1 Flow of the proposed system

4. Implementation And Result

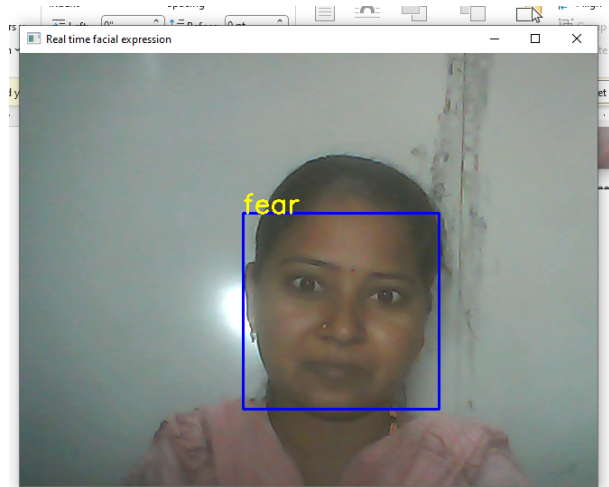


Fig. 2 Fear face detection

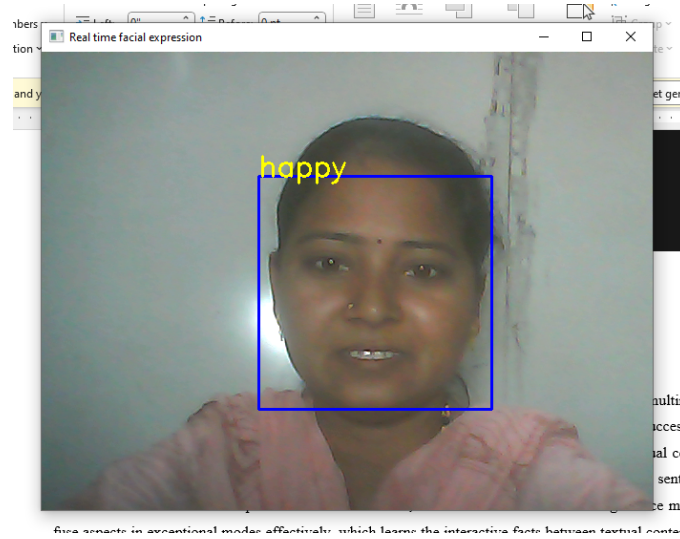


Fig. 3 Happy face detection



Fig.4 Neutral state of sentiment

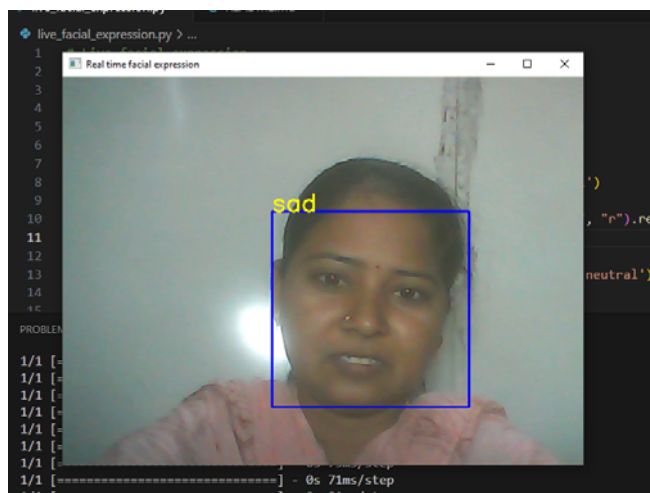


Fig. 5 Sad Face of detection

5. Research Aspects

Multimodal sentiment evaluation of social media is a difficult task. This machine proposes a multimodal sentiment evaluation mannequin based totally on the interest mechanism. This mannequin can successfully dispose of noise interference in the textual statistics of social media and gain extra correct textual content features. Combined with the interest mechanism, the photo aspects that are extra essential to sentiment classification are extracted. In phrases of function fusion, the interest mechanism is brought once more to fuse aspects in exceptional modes effectively, which learns the interactive facts between textual content and images. The mannequin used modal inside facts and modal interplay records to efficaciously acquire the sentimental characteristic illustration of multimodal data, precisely judged the sentimental polarity of users' tweets, higher published users' actual feelings, and helped us apprehend people's attitudes and views closer to sure activities on social media. The experimental outcomes on two open datasets show the feasibility and superiority of our proposed model. In future work, the goal is to enhance the present fashions and techniques and learn about extra modalities, along with audio, video, and so on.

5.1 Mood:

Similar to emotion, moods additionally show off a contagion effect. For example, a depressed individual will frequently make others experience depressed and a completely satisfied man or woman will frequently make others sense happy. Some researcher have proven that even a mere smiling or frowning face, proven so rapidly that the problem is now not aware of seeing the image, can have an effect on a person's temper and because of this bias judgment. From an interface standpoint, the implications for character-based dealers are clear: Moods exhibited via onscreen characters can also immediately switch to the user's mood. Onscreen temper can additionally lead to "perceived contagion" effects: One smiling or frowning face on the display can have an impact on users' perceptions of different faces that they consequently see on the screen, possibly due to the fact of priming.

5.2 Neurological Aspects:

The talent is the most critical supply of emotion. The most frequent way to measure neurological adjustments is the electroencephalogram (EEG). In a blissful state, the human intelligence reveals an alpha rhythm, which can be detected by way of EEG recordings taken via sensors connected to the scalp. Disruption of this sign (alpha blocking) takes place

in response to novelty, complexity, and unexpectedness, as nicely as in the course of emotional exhilaration and anxiety. EEG research have similarly proven that positive/approach-related thoughts lead to higher activation of the left anterior vicinity of the brain, whilst negative/ avoidance-related feelings lead to increased activation of the proper anterior region. Indeed, when one flashes a image to both the left or the proper of the place a man or woman is looking, the viewer can discover a smiling face greater rapidly when it is flashed to the left hemisphere and a frowning face extra shortly when it is flashed to the right hemisphere. Current EEG devices, however, are pretty clumsy and obstructive, rendering them impractical for most HCI applications. Recent advances in magneto resonance imaging (MRI) provide tremendous promise for emotion monitoring, but are presently unrealistic for HCI due to the fact of their expense, complexity, and shape factor.

Table 1 Inputs related to detection of face mask w.r.t. time and accuracy

Sequence	Time (sec)	Accuracy (%)
1	0.90	0.99
2	0.80	0.98
3	0.85	0.99
4	0.96	0.96
5	0.87	0.98
6	0.89	0.99
7	0.79	0.97
8	0.89	0.99
9	0.91	0.95
10	0.92	0.96

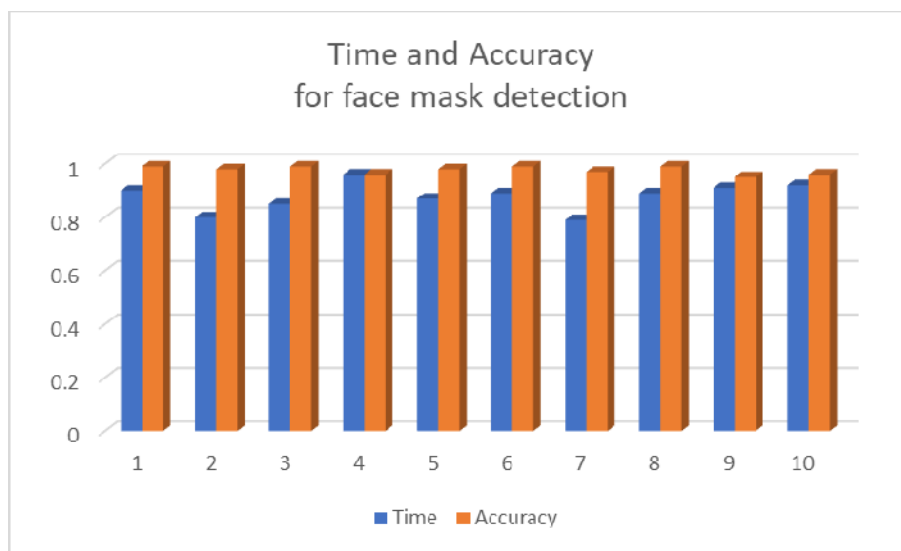


Fig. 6 Graph related to time and accuracy for face mask detection

6. Conclusion

Multi-modal Sentiment Analysis has been getting to know trouble that has been a lookup pastime for current years. Though lot of work is executed until date on sentiment analysis, there are many difficulties to sentiment analyzer on the grounds that Cultural influence, linguistic version and differing contexts make it fairly tough to derive sentiment. The motive at the back of this is unstructured nature of natural language. The foremost difficult components exist in use of different modes; dealing with multi-modality entails the use of a couple of media such as audio and video in addition to textual content to beautify the accuracy of sentiment analyzers. Textual emotional classification is performed on foundation of polarity, depth of lexicons. Audio emotional Classification is performed on groundwork of prosodic features. Video emotional Classification is performed on foundation postures, gestures etc. Infusion, we can combine the outcomes of all these modes; to get greater accuracy. Future lookup should be devoted to these challenges. So we are shifting from uni-modal to multi-modal. At the same time some features related to mood and neurological aspects also been considered to explore the facts of multi modeling sentiment analysis

References

- [1] Efficient Face Recognition System for Operating in Unconstrained Environments Alejandra Sarahi Sanchez-Moreno 1, Jesus Olivares-Mercado 1 , Aldo Hernandez-Suarez 1 , Karina Toscano-Medina 1, Gabriel Sanchez-Perez 1 and Gibran Benitez-Garcia 2,* J. Imaging 2021, 7, 161. <https://doi.org/10.3390/jimaging7090161>
- [2] The Elements of End-to-end Deep Face Recognition: A Survey of Recent Advances Hang Du, Hailin Shi, Dan Zeng, Xiao-Ping Zhang, Tao Mei Accepted for publication in ACM Computing Surveys Computer Vision and Pattern Recognition (cs.CV)
- [3] The Elements of End-to-end Deep Face Recognition: A Survey of Recent Advances HANG DU, HAILIN SHI, DAN ZENG, XIAO-PING ZHANG, TAO MEI, JD AI arXiv:2009.13290v4 [cs.CV] 27 Dec 2021
- [4] (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 9, No. 6, 2018 42 | Page www.ijacsa.thesai.org Study of Face Recognition Techniques: A Survey Madan Lal, Kamlesh Kumar Department of Computer Science Sindh Madressatul Islam University, Karachi, Sindh, Pakistan Rafaqat Hussain Arain, Abdullah Maitlo, Sadaquat Ali Ruk, Hidayatullah Shaikh
- [5] Research on Face Recognition Classification Based on Improved GoogleNet Zhigang Yu, Yunyun Dong, Jihong Cheng, Miaomiao Sun , and Feng Su Hindawi Security and Communication Networks Volume 2022, Article ID 7192306, 6 pages <https://doi.org/10.1155/2022/7192306>
- [6] Boiy, E. Hens, P., Deschacht, K. & Moens, M.F., "Automatic Sentiment Analysis in Online Text", In Proceedings of the Conference on Electronic Publishing(ELPUB-2007).
- [7] J. Wiebe, T. Wilson, and C. Cardie. "Annotating expressions of opinions and emotions in language", Language Resources and Evaluation, 2005.
- [8] Scott Brave and Clifford Nass, Emotion in Human- Computer Interaction. Retrieved from <http://ircm.com.umontreal.ca/dufresne/COM7162/EmotionHumanInteraction.pdf>
- [9] S. V. Bo Pang, Lillian Lee, "Thumbs up? Sentiment classification using machine learning techniques", Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), ACL, pp. 79–86, July 2002.
- [10] Pravesh Kumar Singh and Mohd Shahid Husain "Methodological study of opinion mining and sentiment analysis techniques" International Journal on Soft Computing (IJSC) Vol. 5, No. 1, February 2014

- [11] Gelareh Mohammadi and Alessandro Vinciarelli “Automatic Personality Perception: Prediction of Trait Attribution Based on Prosodic Features”, IEEE transactions on affective computing, vol. 3, no. 3, july-september 2012
- [12] Litman, D.J. and Forbes-Riley, K., “Predicting Student Emotions in Computer-Human Tutoring Dialogues”. In Proc. of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL), July 2004
- [13] Lee C M Narayanan, S.S., “Toward detecting emotions in spoken dialogs”. IEEE Tran. Speech and Audio Processing, Vol. 13 NO. 2, March 2005
- [14] Mozziconacci, S., “Prosody and Emotions”. Int. Conf. on Speech Prosody. 2002 Danielle Lottridge, Mark Chignell, and Michiaki Yasumura,” Identifying Emotion through Implicit and Explicit Measures: Cultural Differences, Cognitive Load, and Immersion”, IEEE transactions on affective computing, vol. 3, no. 2, april-june 2012
- [15] Mohammad Soleymani, Maja Pantic, Thierry Pun,” Multimodal Emotional Recognition in Response to Videos”, IEEE transactions Affective Computing , Vol.3,No.2, April-June 2012
- [16] Gaurav Vasmani and Anuradha Bhatia “A Real Time Approach with Big data- A Survey”, International Journal of Engg Sciences & Research Technology, Vol 3, Issue 9, September 2013
- [17] LouisPhilippe, Morency Rada Mihalcea and Payal Doshi “Towards Multimodal Sentiment Analysis: Harvesting Opinions from the Web”, ICMI’11, November 14–18, 2011, Alicante, Spain.
- [18] Andrea Kleinsmith and Nadia Bianchi-Berthouze ”Affective Body Expression Perception And Recognition:A Survey”, IEEE transactions Affective Computing , Vol.4,No.1, January-March 2013
- [19] Gérard Dray, Michel Plantie , Ali Harb, Pascal Poncelet Mathieu Roche and François Troussel, “Opinion Mining From Blogs” International Journal of Computer Information Systems and Industrial Management Applications – IJCISIM Vol. 1 2009
- [20] Charles B. Ward, Yejin Choi, Steven Skiena, “Empath: A Framework for Evaluating Entity-Level Sentiment Analysis”, IEEE 2011
- [21] Georgios Paltoglou and Michael Thelwall, “Seeing Stars of Valence and Arousal in Blog Posts” IEEE transactions on affective computing, vol. 4, no. 1, January-March 2013
- [22] Shangfei Wang, Zhilei Liu, Zhaoyu Wang, Guobing Wu, Peijia Shen, Shan He, and Xufa Wang “Analyses of a Multimodal Spontaneous Facial Expression Database” IEEE transactions on affective computing, vol. 4, no. 1, January- March 2013
- [23] Ashish Tawari and Mohan Manubhai Trivedi, “Speech Emotion Analysis: Exploring the Role of Context” IEEE transactions on multimedia, vol. 12, no. 6, October 2010
- [24] Tal Sobol-Shikler and Peter Robinson, “Classification of Complex Information: Inference of Co-Occurring Affective States from Their Expressions in Speech”, IEEE transactions on pattern analysis and machine intelligence, vol. 32,no.7,July 2010
- [25] Felix Weninger, Jarek Krajewski, Anton Batliner, and Bjorn Schuller, “The Voice of Leadership:Models and Performances of Automatic Analysis in Online Speeches”, IEEE transactions on affective computing, vol. 3, no. 4, October-December.2012