# Implementation of Visual Sentiment Analysis on Flickr Images

### Harshala Bhoir[1], Dr. K. Jayamalini[2]

[1]Department of Computer Engineering, Shree.L.R.Tiwari College of Engineering ,Mira road(E),Thane, Maharashtra, India

[2]Assistant Professor Department of Computer Engineering, Shree.L.R.Tiwari College Of Engineering, Mira road(E),Thane, Maharashtra, India

## ABSTRACT

Visual sentiment analysis is the way to automatically recognize positive and negative emotions from images, videos, graphics and stickers. To estimate the polarity of the sentiment evoked by images in terms of positive or negative sentiment, most of the state of the art works exploit the text associated with a social post provided by the user. However, such textual data is typically noisy due to the subjectivity of the user, which usually includes text useful to maximize the diffusion of the social post. This System will extract three views: visual view, subjective text view and objective text view of Flickr images and will give sentiment polarity positive, negative or neutral based on the hypothesis table. Subjective text view gives sentiment polarity using VADER (Valence Aware Dictionary and sEntiment Reasoner) and objective text view gives sentiment polarity with three convolution neural network models. This system implements VGG-16, Inception-V3 and ResNet-50 convolution neural networks with pre pre-trained ImageNet dataset. The text extracted through these three convolution networks is given to VADER as input to find sentiment polarity. This system implements visual view using a bag of visual word model with BRISK (Binary Robust Invariant Scalable Key points) descriptor. System has a training dataset of 30000 positive, negative and neutral images. All the three views' sentiment polarity is compared. The final sentiment polarity is calculated as positive if two or more views gives positive sentiment polarity, as negative if two or more views gives negative sentiment polarity and as neutral if two or more views gives neutral sentiment polarity. If all three views give unique polarity then the polarity of the objective text view is given as output sentiment polarity.

Keywords: Sentiment analysis, CNN, ResNet-50, Inception-V3, VGG1-16,Bag of visual words, Vader, Feature Extraction, subjective text view, objective text view, BRISK, Imagenet, Keras, Tenserflow

## I. INTRODUCTION

### 1.1 Visual Sentiment Analysis

Visual sentiment analysis is the way to automatically recognize positive and negative emotions from images, videos, graphics and stickers. Sentiment analysis is the automated process of understanding an opinion about a given subject from written or spoken language. In a world where we generate 2.5 quintillion bytes of data every day, sentiment analysis has become a key tool for making sense of that data [7].With the popularity of social networks and mobile devices, there is a huge volume of images and videos captured by users to record all kinds of activities in their lives every day and everywhere. For example, people may share their travel experiences, their opinions towards some events and so on. Automatically analysing the sentiment from these multimedia contents is demanded by many practical applications, such as smart advertising, targeted marketing and political voting forecasts. Compared with text based sentiment analysis which infers emotional signals from short textual description, visual contents, such as colour contrast and tone could provide more vivid clues to reveal the sentiment behind.

Figure 1 shows examples of images with positive, negative or neutral sentiment. Apparently, the images in the upper row manifest positive sentiment, while those in the middle row deliver negative sentiment and images in the lower row have neutral sentiment.



Figure 1: Example of images with positive, negative or neutral sentiment

### 1.2 Problem statement

This system does visual sentiment analysis on live Flickr images by extracting three views of the input image: two text views and one visual view. First text view is subjective in which the title provided with the Flickr image is taken as input and fed to VADER for sentiment analysis. Second text view is objective where VGG-16, Inception-V3 and ResNet-50 CNNs are applied on Flickr images to extract text related to the image rather than reading the title and gives sentiment using VADER. This system also generates a visual view using a bag of visual words image classifier with BRISK descriptor to get sentiment. After implementation of sentiment analysis on text and visual view final output was obtained using table 1 hypothesis table.

### 1.3 Project Objective

To do visual sentiment analysis on Flickr images by extracting and employing an Objective Text description of images automatically extracted from the visual content rather than the classic Subjective Text provided by the users and visual view of image based on hypothesis table. Following are objectives of project:-

- To do visual sentiment analysis on Flickr image by extracting titles provided from the user.
- To do visual sentiment analysis on Flickr image data using VGG-16,Inception-V3 and ResNet-50 CNNs
- To do visual sentiment analysis on Flickr image by extracting visual view using BOVW with BRISK descriptor.
- To get final sentiment polarity of Flicker image using hypothesis table by considering polarity of text view and visual view.

### 1.4 Project Idea

Social media users continuously post images together with their opinions and share their emotions. This trend has supported the growth of new application areas, such as semantic based image selection from

crowd-sourced collections [1], Social Event Analysis [1] and Sentiment Analysis on Visual Contents [9]. Visual Sentiment Analysis aims to infer the sentiment evoked by images in terms of positive or negative polarity. Early methods in this field focused only on visual features or have employed text to define a sentiment ground truth. More recent approaches combine visual and text features by exploiting well-known semantic and sentiment lexicons [1]

In the proposed system the text associated with images is typically obtained by considering the meta-data provided by the user (e.g., image title, tags and description). It also describes images in an "objective" way by using scene understanding methods [1].Automatic extraction of text by the system from image is called objective text. The "objective" emphasizes the fact that it is different to the" subjective" text written by the user for an image of a post.The system will use VGG-16, Inception-V3 and ResNet-50 CNNs architectures to extract text from the image. Also the system uses a Bag of visual words image descriptor and BRISK descriptor for visual view of image. The proposed system will extract three views of given social media image i.e. Visual view, subjective text view and objective text view will give sentiment polarity based on the given hypothesis table shown in table 1 using rules.

## II.  Literature Review

Several papers investigated the problem of joint modelling the representation of images and associated text or tags for different tasks, such as image retrieval, social images understanding, image annotation and visual sentiment analysis. In [1] the authors presented different learning architectures for sentiment analysis of social posts containing short text messages and an image (i.e., Tweets). They exploited a representation learning architecture that combines the input text with the polarity ground truth.

The approach proposed in [10], combines visual features with text-based features extracted from the text subjectively associated by the users to images (i.e., descriptions and tags). To represent contents for sentiment analysis estimation, the authors proposed three different type of features extracted considering pairs of images and the related subjective texts: a visual feature defined by combining different visual descriptors usually used for visual classification ,a feature obtained by using the traditional Bag of Words approach on the subjective text, and a sentiment feature obtained by selecting the words of the subjective text whose sentiment scores (positive or negative) reported in SentiWordNet  are larger than a threshold, and applying the Bag of Words on this restricted vocabulary. The considered features are exploited to define embedding space in which the correlation among the projected features is maximized. Then a sentiment classifier is trained on the features projected in the embedding space.

In [1] the author  extracts and employs an Objective Text description of images automatically extracted from the visual content rather than the classic Subjective Text provided by the users. The proposed method defines a multimodal embedding space based on the contribution of both visual and textual features. The sentiment polarity is then inferred by a supervised Support Vector Machine trained on the representations of the obtained embedding space. Experiments performed on a representative dataset of 47235 labelled samples demonstrate that the exploitation of the proposed Objective Text helps to outperform state-of-the-art for sentiment polarity estimation.

In [8] paper, an image dataset with sentiment tags is built for training .Authors conduct experiments by training 15000 scene images on three different CNNs models Inception V3, ResNet and VGG Net and providing that deep learning can perform rather well on sentiment prediction tasks.

## III. Proposed system

Section 3.1 describes features extraction process, Section 3.2 describes proposed System diagram and section 3.3 describes expected results.

### 3.1 Feature Extraction

The proposed approach exploits one visual view and two textual views based on the objective text extracted from the images and Subjective Text provided by the user with the image. The following subsections details the feature extraction process. Figure 2 shows feature extraction used in the proposed system.

### 3.1.1 Visual View

System uses a Bag of visual words image classifier with BRISK (Binary Robust Invariant Scalable Key points) descriptor for visual representation. The BOVW uses a training set of 30000 images with positive, negative and neutral labels.

### 3.1.2 Text View

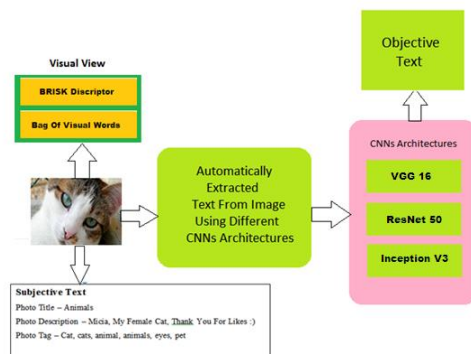There are two textual views based on Subjective Text and objective text extracted from the images.



Figure 2: Feature Extraction of System

### Subjective Text view

This view reflects the subjective text information provided by the users such as photo title, description and tags. It consists of textual features, which will extract from text associated with image This text is send as a input to VADER (Valence Aware Dictionary and sEntiment Reasoner) .It is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media. VADER uses a combination of sentiment lexicon is a list of lexical features (e.g., words) which are generally labelled according to their semantic orientation as either positive or negative. VADER a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.

### Objective Text View

Objective Text will be obtained through three deep learning CNNs architectures VGG-16, Inception-V3 and ResNet-50. As shown in Figure 3.2, each architecture will provide a description, in some sense objective, of the input image from a different point of view, as each architecture has been trained for a different task. This will allow to obtain a wide objective description of the image which takes into account different semantic aspects of the visual content. Redundant terms are not a drawback for the proposed approach; indeed the presence of more occurrences of similar or related terms enhances the weight of these correct terms in the representation extracted by the proposed system, and reduces the effect of noisy results. For these reasons the system will find unique words of three CNNs output text and send this as input to VADER to find sentiment polarity of image.

This section presents a visual sentiment analysis method that uses visual view and text views. Figure 3 shows the system diagram.
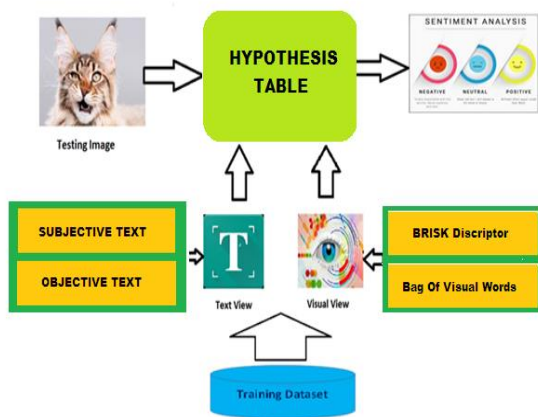
## 3.2 Proposed System Diagram



Figure 3: System Diagram

As shown, System will first extract features from each view and then do visual sentiment analysis.

The system has visual views with BRISK descriptor and a Bag of visual words image descriptor visual representation. The system has two text views such as subjective text view (provided by user), Objective text view (Extracted from image using different CNNs architectures i.e. VGG-16, Inception-V3 and ResNet-50)

After feature extraction from all views, visual sentiment analysis is done to give sentiment polarity "positive" or "negative" or "Neutral" based on the hypothesis table as shown in table I using four rules describe in section 3.3.

## 3.3 Expected Result

The proposed system will extract three views of a given social media image i.e. visual view, subjective text view and objective text view and will give sentiment polarity based on the given hypothesis table shown in table 1 using following rules.

**Rule 1 :** If any two or more views among three views have positive polarity then the proposed system will give output as positive sentiment polarity as shown in serial numbers 1,2,3,4,7,10,19.

**Rule 2 :** If any two or more views among three views have negative polarity then the proposed system will give output as negative sentiment polarity as shown in serial numbers 5,11,13,14,15,17,23.

**Rule 3 :** If any two or more views among three views have neutral polarity then the proposed system will give output as neutral sentiment polarity as shown in serial numbers 9,18,21,24,25,26,27.

**Rule 4 :** If all three views of image have unique polarity i.e. one positive ,one negative and one neutral polarity then the system will consider objective text view polarity as output polarity as shown in highlighted serial numbers 6,8,12,16,20,22.

TABLE 1 HYPOTHESIS TABLE FOR SENTIMENT POLARITY

| SR. NO | Sentiment Polarity | | | Proposed System sentiment Polarity |
|---|---|---|---|---|
| | Subjective Text View | Objective Text View | Visual View | |
| 1 | POSITIVE | POSITIVE | POSITIVE | POSITIVE |
| 2 | POSITIVE | POSITIVE | NEGATIVE | POSITIVE |
| 3 | POSITIVE | POSITIVE | NEUTRAL | POSITIVE |
| 4 | POSITIVE | NEGATIVE | POSITIVE | POSITIVE |
| 5 | POSITIVE | NEGATIVE | NEGATIVE | NEGATIVE |
| 6 | POSITIVE | NEGATIVE | NEUTRAL | NEGATIVE |
| 7 | POSITIVE | NEUTRAL | POSITIVE | POSITIVE |
| 8 | POSITIVE | NEUTRAL | NEGATIVE | NEUTRAL |
| 9 | POSITIVE | NEUTRAL | NEUTRAL | NEUTRAL |
| 10 | NEGATIVE | POSITIVE | POSITIVE | POSITIVE |

| 11 | NEGATIVE | POSITIVE | NEGATIVE | NEGATIVE |
| 12 | NEGATIVE | POSITIVE | NEUTRAL | POSITIVE |
| 13 | NEGATIVE | NEGATIVE | POSITIVE | NEGATIVE |
| 14 | NEGATIVE | NEGATIVE | NEGATIVE | NEGATIVE |
| 15 | NEGATIVE | NEGATIVE | NEUTRAL | NEGATIVE |
| 16 | NEGATIVE | NEUTRAL | POSITIVE | NEUTRAL |
| 17 | NEGATIVE | NEUTRAL | NEGATIVE | NEGATIVE |
| 18 | NEGATIVE | NEUTRAL | NEUTRAL | NEUTRAL |
| 19 | NEUTRAL | POSITIVE | POSITIVE | POSITIVE |
| 20 | NEUTRAL | POSITIVE | NEGATIVE | POSITIVE |
| 21 | NEUTRAL | POSITIVE | NEUTRAL | NEUTRAL |
| 22 | NEUTRAL | NEGATIVE | POSITIVE | NEGATIVE |
| 23 | NEUTRAL | NEGATIVE | NEGATIVE | NEGATIVE |
| 24 | NEUTRAL | NEGATIVE | NEUTRAL | NEUTRAL |
| 25 | NEUTRAL | NEUTRAL | POSITIVE | NEUTRAL |
| 26 | NEUTRAL | NEUTRAL | NEGATIVE | NEUTRAL |
| 27 | NEUTRAL | NEUTRAL | NEUTRAL | NEUTRAL |

## IV. Implementation

This section has detailed methods that are used to implement the project. This project takes live input images from the Flickr website through the Flickr API and gives sentiment polarity based on subjective text, objective text and visual view.

Subjective text view is extracted from the Flickr images by reading the title of the image provided by the user. This view reads the title and gives the title as an input to VADER ((Valence Aware Dictionary and sEntiment Reasoner) to get sentiment polarity. Since Vader is optimized for social media data also It is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media. I selected VADER instead of TextBlob for text view sentiment analysis. Both the TextBlob and VADER have 56% accuracy.

Objective text view is extracted from the Flickr images directly from three convolution neural networks. This application implemented VGG16, Inception V3 and ResNet 50.All these CNNs are pre trained on the ImageNet database. ImageNet is an image database organized according to the WordNet hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. CNN gives the first five predicted words related to image. All the words from three CNNs are compared and unique words are fed to VADER as input to get sentiment polarity.

### Accuracy of VGG–16, ResNet-50 and Inception-V3:-

Keras Applications are deep learning models that are made available alongside pre-trained weights. These models can be used for prediction, feature extraction, and fine-tuning .Table II shows top-1 and top-5 accuracy refers to the model's performance on the ImageNet validation dataset.

| Model | Top-1 Accuracy | Top-5 Accuracy | Parameters | Year |
|---|---|---|---|---|
| VGG-16 | 71.3% | 90.1% | 138,357,544 | 2014 |
| ResNet-50 | 74.9% | 92.1% | 25,636,712 | 2015 |
| Inception-V3 | 77.9% | 93.7% | 23,851,784 | 2015 |

TABLE II TOP-1 AND TOP-5 ACCURACY OF VGG-16, RESNET-50 AND INCEPTION-V3

In visual view Bag of visual model is used to classify the image as positive, negative and neutral. In the visual view I used BRISK ( Binary Robust Invariant Scalable Key points) instead of SIFT (scale-invariant feature transform) descriptor. BRISK relies on an easily configurable circular sampling pattern from which it computes brightness comparisons to form a binary descriptor string. The unique properties of BRISK can be useful for a wide spectrum of applications, in particular for tasks with hard real-time constraints or limited computation power: BRISK finally offers the quality of high-end features in such time-demanding applications [15]. BRISK detects more features than SIFT.

## 4.1 Data Collection

**Flicker Images** - To fetch Flickr image data, I used flickrapi module for python .I used access key and secret key to authorize Flickr account .To fetch data for users input keyword (text) I used photo.search() method.

**Training Dataset** - This data is used to train the classifier. To collect this data I used the NLTK library of python. Sample of this training set is shown in figure 4, figure 5 and figure 6.Training set has 30000 training images.
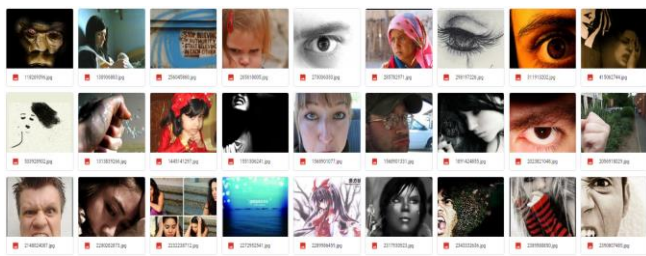


Figure 4 Sample Negative images in training dataset



Figure 5 Sample Positive images in training dataset



Figure 6 Sample Neutral images in training dataset

## V. Result

In this section I am going to show various outputs that I got in project implementation.

Figure 7 shows 25 sample Input images fetched through Flickr API. I will extract features from input images. System extracted subjective text view sentiment polarity using VADER as shown in figure 8, objective text view sentiment polarity using output of top 5 words extracted through VGG-16 ,Inception-V3 and ResNet-50 CNN as shown in figure 9.System also extracted visual view sentiment polarity using BOVW with BRISK image descriptor as shown in figure 11.
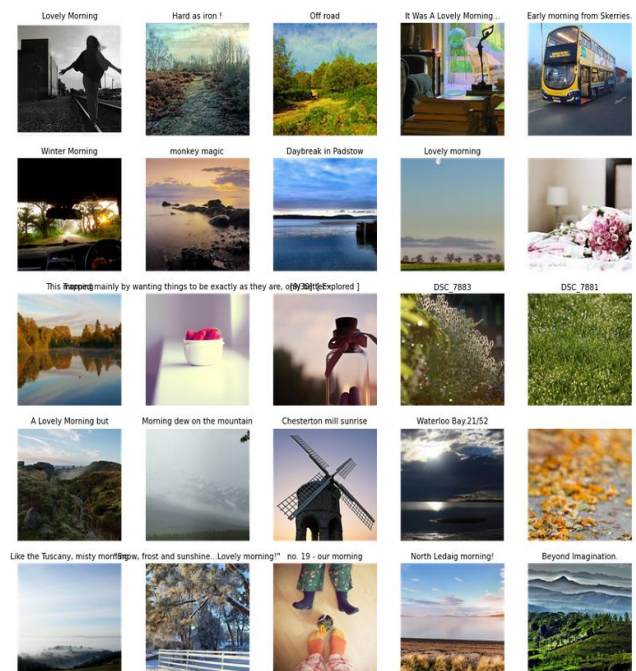


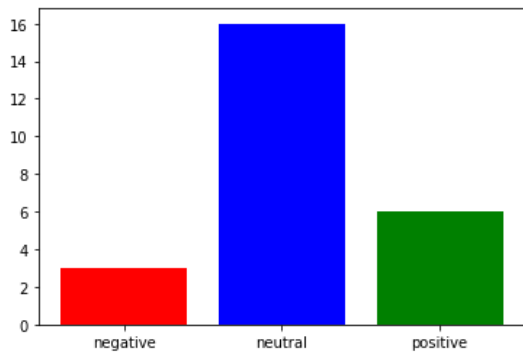Figure 7 Input images fetched through Flickr API

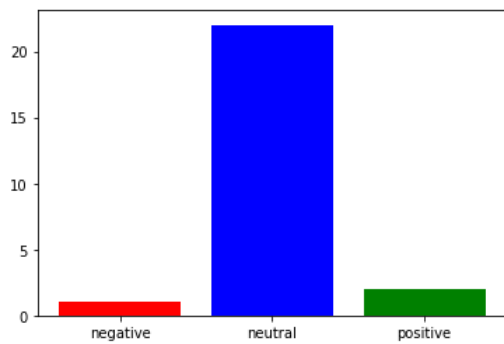Figure 8 Visual Sentiment analysis of subjective text view



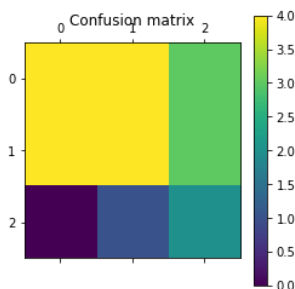Figure 9 Visual Sentiment analysis of objective text view



Figure 10 Confusion matrix of visual view

Figure 10 shows the confusion matrix of 25 Flickr images. These 25 images are classified using a bag of visual words classifier with BRISK descriptor. For the visual view I used the training data set bovw-2 .It has approximately 30000 positive, negative and neutral images. Accuracy measure is one of the most important steps in Machine learning algorithms .A Confusion matrix is basically how many test cases were correctly classified. Hence, a confusion matrix is used to determine the accuracy of classification.
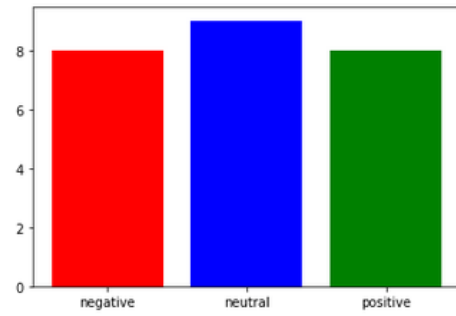


Figure 11 Visual Sentiment analysis visual view

Final result is given as per rules in hypothesis table shown in Table I in section 3.3. Figure 12 shows the bar graph of visual sentiment analysis of subjective view with positive, negative and neutral sentiment polarity of 25 input images taken through Flickr API.
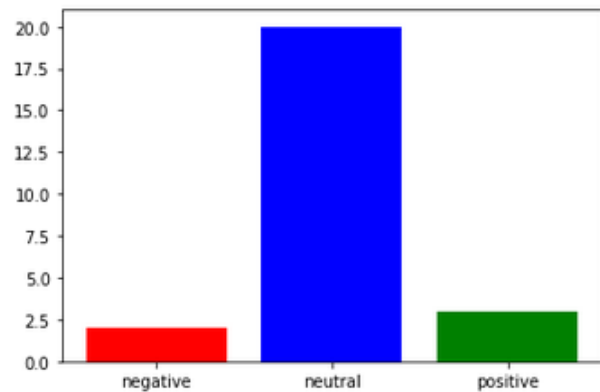


Figure 12 Final output of Visual Sentiment analysis with text view and visual view

## VI. Conclusion

This system addresses the challenge of image sentiment polarity estimation by proposing a novel source of text for this task. The aim is to deal with the issue related to the text provided by users which is commonly used in most of the previous works. It explained a study in which Objective Text extracted considering the visual content of images is compared with respect to the Subjective Text provided by users. This system first identified several drawbacks brought by the Subjective Text due its intrinsic nature, and then it demonstrated experimentally that the exploitation of Objective Text associated with images

provides better results than the use of the Subjective Text provided by the user. The Objective Text exploited by the proposed approach will not present the highlighted limitations and it will automatically extract from the image. The expected result will support the use of Objective text automatically extracted from images for the task of Visual Sentiment Analysis in lieu of the Subjective Text provided by users and given sentiment polarity based on hypothesis table 3.1. Subjective text view gave sentiment polarity using VADER and objective text view gave sentiment polarity with three convolution neural network models. This system implemented VGG-16, Inception-V3 and ResNet-50 convolution neural networks. The text extracted through these three convolution networks fed to VADER as input to find sentiment polarity. In this system visual view is implemented with a bag of visual word models using BRISK descriptor having a training set of approximately 30000 images. All the three view sentiment polarity is compared. The final sentiment polarity is calculated as positive if two or more views gives positive sentiment polarity, as negative if two or more views gives negative sentiment polarity and as neutral if two or more views gives neutral sentiment polarity. If all three views give unique polarity then the polarity of the objective text view is given as output sentiment polarity.

## VII. REFERENCES

[1]. Alessandro Ortis Giovanni M. Farinella,Giovanni Torrisi,Sebastiano Battiato Visual Sentiment Analysis Based on Objective Text Journal]. – Catania, Italy : IEEE, 2018. – Vols. 978-1-5386-7021-7/18/.

[2]. B. Zhou A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva Learning deep features for scene recognition using places database Journal]. – 3Universitat Oberta de Catalunya : s.n.]. – Vols. Advances in Neural Information Processing Systems, 2014, pp. 487–495.

[3]. Bertini1 Claudio Baecchi1 ·Tiberio Uricchio1 · Marco A multimodal feature learning approach for sentiment Journal]. – New York : Springer, 2015

[4]. C. Szegedy W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, Going deeper with convolutions Journal]. – s.l.] : IEEE, 2015. – Vol. In proceedings of the IEEE Conference on Computer Vision and Pattern.

[5]. Eunjeong Ko Chanhee Yoon,Eun Yi Kim Discovering Visual Features for Recognizing User's Journal]. – Konkuk University,South Korea : IEEE, 2016. – Vols. 978-1-4673-8796-5/16.

[6]. Fei-Fei A. Karpathy and L. Deep visual-semantic alignments for generating image descriptions Journal]. – s.l.] : IEEE. – Vols. JOURNAL OF LATEX CLASS FILES, VOL. 14, NO. 8, AUGUST 2015.

[7]. Vikrant Waghmare, Mahesh Pimpalkar, Prof. Vaishali Londhe, Automated analysis techniques to extract sentiments and opinions conveyed in the user comments on social mediaJournal] JETIR, Yadavrao Tasgaonkar Institute of Engineering & Technology University of Mumbai December 2018, Volume 5, Issue 12

[8]. Junfeng Yao Yao Yu and Xiaoling Xue Sentiment Prediction In Scene Images Via Convolution Neural Networks Journal]. – Beijing,China : IEEE, 2016. – Vols. 978-1-5090-4423-8/16.

[9]. Kaikai Songa Ting Yaob, Qiang Linga,∗, Tao Mei Boosting Image Sentiment Analysis with Visual Attention Journal]. – China : ELSEWHERE, 2018.

[10]. Marie Katsurai Shin'ichi Satoh IMAGE SENTIMENT ANALYSIS USING LATENT CORRELATIONS AMONG VISUAL, Journal]. – Tokyo, Japan : IEEE, 2016. – Vols. 978-1-4799-9988-0/16.

[11]. Varshney Mayank Amencherla and Lav R. Color-Based Visual Sentiment for Social Journal]. - Urbana-Champaign : IEEE, 2017. - Vols. 978-1-5090-6026-9/17.

[12]. Vincent Feng An Overview of ResNet and its Variants ,Jul 16,2017 Available: https://towardsdatascience.com/an-overview-of-resnet-and-its-variants-5281e2f56035 , Last visit 4/10/21

[13]. Muneeb ul Hassan ,VGG16 – Convolutional Network for Classification and Detection,20 November 2018 Available : https://neurohive.io/en/popular-networks/vgg16/ , Last visit: 4/10/21

[14]. Aqeel Anwar, Difference between AlexNet, VGGNet, ResNet, and Inception,Jun 7 2021 Available: https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96 , Last visit: 9/10/21

[15]. Stefan Leutenegger et.al BRISK: Binary Robust Invariant Scalable Keypoints- IEEE,2011 –Vols 978-1-4577-1102-2/11

[16]. Valeria Maeda-Gutiérrez et.al Comparison of Convolutional Neural Network Architectures for Classification of Tomato Plant Diseases Appl. Sci. 2020, 10, 1245; doi:10.3390/app10041245

[17]. Michele Compri MULTI-LABEL REMOTE SENSING IMAGE RETRIEVAL BASED ON DEEP FEATURES 2016,universita DEGLI STUDI DI TRENTO

[18]. Hussain Mujtaba , Introduction to Resnet or Residual Network .Sep 28,2020 Available: https://www.mygreatlearning.com/blog/resnet/#sh1 , Last visit: 9/10/21

[19]. ResNet,AlexNet,VGG Net ,Inception :Understanding Various Architectures of covolution neural networks,Available : https://cv-tricks.com/cnn/understand-resnet-alexnet-vgg-inception/ , Last visit: 9/10/21

[20]. Zharfan Zahisham et.al Food Recognition with ResNet-50 ,IEEE, 2020 – Vols. 978-1-7281-6946-0/20

## Cite this article as :

Harshala Bhoir, Dr. K. Jayamalini, "Implementation of Visual Sentiment Analysis on Flickr Images", International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT), ISSN : 2456-3307, Volume 7 Issue 5, pp. , September-October 2021. Available at doi : https://doi.org/10.32628/CSEIT217533
Journal URL : https://ijsrcseit.com/CSEIT217533